# MODEL REDUCTION OF DESCRIPTOR SYSTEMS BY INTERPOLATORY PROJECTION METHODS

SERKAN GUGERCIN[*], TATJANA STYKEL[†], AND SARAH WYATT[‡]

**Abstract.** In this paper, we investigate interpolatory projection framework for model reduction of descriptor systems. With a simple numerical example, we first illustrate that employing subspace conditions from the standard state space settings to descriptor systems generically leads to unbounded $\mathcal{H}_2$ or $\mathcal{H}_\infty$ errors due to the mismatch of the polynomial parts of the full and reduced-order transfer functions. We then develop modified interpolatory subspace conditions based on the deflating subspaces that guarantee a bounded error. For the special cases of index-1 and index-2 descriptor systems, we also show how to avoid computing these deflating subspaces explicitly while still enforcing interpolation. The question of how to choose interpolation points optimally naturally arises as in the standard state space setting. We answer this question in the framework of the $\mathcal{H}_2$-norm by extending the Iterative Rational Krylov Algorithm (IRKA) to descriptor systems. Several numerical examples are used to illustrate the theoretical discussion.

**Key words.** interpolatory model reduction, differential algebraic equations, $\mathcal{H}_2$ approximation

**AMS subject classifications.** 41A05, 93A15, 93C05, 37M99

**1. Introduction.** We discuss interpolatory model reduction of differential-algebraic equations (DAEs), or descriptor systems, given by

$$\begin{aligned}
\mathbf{E}\,\dot{\mathbf{x}}(t) &= \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t), \\
\mathbf{y}(t) &= \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{u}(t),
\end{aligned} \tag{1.1}$$

where $\mathbf{x}(t) \in \mathbb{R}^n$, $\mathbf{u}(t) \in \mathbb{R}^m$ and $\mathbf{y}(t) \in \mathbb{R}^p$ are the states, inputs and outputs, respectively, $\mathbf{E} \in \mathbb{R}^{n \times n}$ is a *singular* matrix, $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\mathbf{B} \in \mathbb{R}^{n \times m}$, $\mathbf{C} \in \mathbb{R}^{p \times n}$, and $\mathbf{D} \in \mathbb{R}^{p \times m}$. Taking the Laplace transformation of system (1.1) with zero initial condition $\mathbf{x}(0) = \mathbf{0}$, we obtain $\widehat{\mathbf{y}}(s) = \mathbf{G}(s)\widehat{\mathbf{u}}(s)$, where $\widehat{\mathbf{u}}(s)$ and $\widehat{\mathbf{y}}(s)$ denote the Laplace transforms of $\mathbf{u}(t)$ and $\mathbf{y}(t)$, respectively, and $\mathbf{G}(s) = \mathbf{C}(s\mathbf{E} - \mathbf{A})^{-1}\mathbf{B} + \mathbf{D}$ is a transfer function of (1.1). By following the standard abuse of notation, we will denote both the dynamical system and its transfer function by $\mathbf{G}$.

Systems of the form (1.1) with extremely large state space dimension $n$ arise in various applications such as electrical circuit simulations, multibody dynamics, or semidiscretized partial differential equations. Simulation and control in these large-scale settings is a huge computational burden. Efficient model utilization becomes crucial where model reduction offers a remedy. The goal of model reduction is to replace the original dynamics in (1.1) by a model of the same form but with much smaller state space dimension such that this reduced model is a high fidelity approximation to the original one. Hence, we seek a reduced-order model

$$\begin{aligned}
\widetilde{\mathbf{E}}\,\dot{\widetilde{\mathbf{x}}}(t) &= \widetilde{\mathbf{A}}\widetilde{\mathbf{x}}(t) + \widetilde{\mathbf{B}}\mathbf{u}(t), \\
\widetilde{\mathbf{y}}(t) &= \widetilde{\mathbf{C}}\widetilde{\mathbf{x}}(t) + \widetilde{\mathbf{D}}\mathbf{u}(t),
\end{aligned} \tag{1.2}$$

where $\widetilde{\mathbf{E}}, \widetilde{\mathbf{A}} \in \mathbb{R}^{r \times r}$, $\widetilde{\mathbf{B}} \in \mathbb{R}^{r \times m}$, $\widetilde{\mathbf{C}} \in \mathbb{R}^{p \times r}$, and $\widetilde{\mathbf{D}} \in \mathbb{R}^{p \times m}$ such that $r \ll n$, and the error $\mathbf{y} - \widetilde{\mathbf{y}}$ is small with respect to a specific norm over a wide range of inputs $\mathbf{u}(t)$

---
[*]Serkan Gugercin is with the Department of Mathematics, Virginia Tech., Blacksburg, VA, 24061-0123, USA, e-mail: `gugercin@math.vt.edu`.

[†]Tatjana Stykel is with Institut für Mathematik, Universität Augsburg, Universitätsstraße 14, 86159 Augsburg, Germany, e-mail: `stykel@math.uni-augsburg.de`.

[‡]Sarah Wyatt is with the Department of Mathematics, Indian River State College, Fort Pierce, FL, 34981, USA, e-mail: `swyatt@irsc.edu`.

with bounded energy. In the frequency domain, this means that the transfer function of (1.2) given by $\widetilde{\mathbf{G}}(s) = \widetilde{\mathbf{C}}(s\widetilde{\mathbf{E}} - \widetilde{\mathbf{A}})^{-1}\widetilde{\mathbf{B}} + \widetilde{\mathbf{D}}$ approximates $\mathbf{G}(s)$ well, i.e., the error $\mathbf{G}(s) - \widetilde{\mathbf{G}}(s)$ is small in a certain system norm.

The reduced-order model (1.2) can be obtained via projection as follows. We first construct two $n \times r$ matrices $\mathbf{V}$ and $\mathbf{W}$, approximate the full-order state $\mathbf{x}(t)$ by $\mathbf{V}\widetilde{\mathbf{x}}(t)$, and then enforce the Petrov-Galerkin condition

$$\mathbf{W}^T \left( \mathbf{E}\mathbf{V}\dot{\widetilde{\mathbf{x}}}(t) - \mathbf{A}\mathbf{V}\widetilde{\mathbf{x}}(t) - \mathbf{B}\,\mathbf{u}(t) \right) = \mathbf{0}, \qquad \widetilde{\mathbf{y}}(t) = \mathbf{C}\mathbf{V}\widetilde{\mathbf{x}}(t) + \mathbf{D}\mathbf{u}(t).$$

As a result, we obtain the reduced-order model (1.2) with the system matrices

$$
\begin{aligned}
\widetilde{\mathbf{E}} &= \mathbf{W}^T\mathbf{E}\mathbf{V}, & \widetilde{\mathbf{A}} &= \mathbf{W}^T\mathbf{A}\mathbf{V}, \\
\widetilde{\mathbf{B}} &= \mathbf{W}^T\mathbf{B}, & \widetilde{\mathbf{C}} &= \mathbf{C}\mathbf{V}, & \widetilde{\mathbf{D}} &= \mathbf{D}.
\end{aligned}
\tag{1.3}
$$

The projection matrices $\mathbf{V}$ and $\mathbf{W}$ determine the subspaces of interest and can be computed in many different ways.

In this paper, we consider projection-based interpolatory model reduction methods, where the choice of $\mathbf{V}$ and $\mathbf{W}$ enforces certain tangential interpolation of the original transfer function. These methods will be presented in Section 2 in more detail. Projection-based interpolation with multiple interpolation points was initially proposed by Skelton *et. al.* in [7, 31, 32]. Grimme [10] has later developed a numerically efficient framework using the rational Krylov subspace method of Ruhe [24]. The tangential rational interpolation framework, we will be using here, is due to a recent work by Gallivan *et al.* [9].

Unfortunately, it is often assumed that extending interpolatory model reduction from standard state space systems with $\mathbf{E} = \mathbf{I}$ to descriptor systems with singular $\mathbf{E}$ is as simple as replacing $\mathbf{I}$ by $\mathbf{E}$. In Section 2, we present an example showing that this naive approach may lead to a poor approximation with an unbounded error $\mathbf{G}(s) - \widetilde{\mathbf{G}}(s)$ although the classical interpolatory subspace conditions are satisfied. In Section 3, we modify these conditions in order to enforce bounded error. The theoretical result will take advantage of the spectral projectors. Then using the new subspace conditions, we extend in Section 4 the optimal $\mathcal{H}_2$ model reduction method of [15] to descriptor systems. Sections 3 and 4 make explicit usage of deflating subspaces which could be numerically demanding for general problems. Thus, for the special cases of index-1 and index-2 descriptor systems, we show in Sections 5 and 6, respectively, how to apply interpolatory model reduction without explicitly computing the deflating subspaces. Theoretical discussion will be supported by several numerical examples. In particular, in Section 5.2, we present an example, where the balanced truncation approach [26] is prone to failing due to problems solving the generalized Lyapunov equations, while the (optimal) interpolatory model reduction can be effectively applied.

**2. Model reduction by tangential rational interpolation.** The goal of model reduction by tangential interpolation is to construct a reduced-order model (1.2) such that its transfer function $\widetilde{\mathbf{G}}(s)$ interpolates the original one, $\mathbf{G}(s)$, at selected points in the complex plane along selected directions. We will use the notation of [1] to define this problem more precisely: Given $\mathbf{G}(s) = \mathbf{C}(s\mathbf{E} - \mathbf{A})^{-1}\mathbf{B} + \mathbf{D}$, the left interpolation points $\{\mu_i\}_{i=1}^q$, $\mu_i \in \mathbb{C}$, together with the left tangential directions $\{\mathsf{c}_i\}_{i=1}^q$, $\mathsf{c}_i \in \mathbb{C}^p$, and the right interpolation points $\{\sigma_j\}_{j=1}^r$, $\sigma_j \in \mathbb{C}$, together with the

right tangential directions $\{b_j\}_{j=1}^r$, $b_j \in \mathbb{C}^m$, we seek to find a reduced-order model $\widetilde{\mathbf{G}}(s) = \widetilde{\mathbf{C}}(s\widetilde{\mathbf{E}} - \widetilde{\mathbf{A}})^{-1}\widetilde{\mathbf{B}} + \widetilde{\mathbf{D}}$ that is a tangential interpolant to $\mathbf{G}(s)$, i.e.,

$$
\begin{aligned}
c_i^T \mathbf{G}(\mu_i) &= c_i^T \widetilde{\mathbf{G}}(\mu_i), & i &= 1, \ldots, q, \\
\mathbf{G}(\sigma_j) b_j &= \widetilde{\mathbf{G}}(\sigma_j) b_j, & j &= 1, \ldots, r.
\end{aligned}
\tag{2.1}
$$

Through out the paper, we will assume $q = r$, meaning that the same number of left and right interpolation points are used. In addition to interpolating $\mathbf{G}(s)$, one might ask for matching the higher-order derivatives of $\mathbf{G}(s)$ along the tangential directions as well. This scenario will also be handled.

By combining the projection-based reduced-order modeling technique with the interpolation framework, we want to find the $n \times r$ matrices $\mathbf{W}$ and $\mathbf{V}$ such that the reduced-order model (1.2), (1.3) satisfies the tangential interpolation conditions (2.1). This approach is called projection-based interpolatory model reduction. How to enforce the interpolation conditions via projection is shown in the following theorem, where the $\ell$-th derivative of $\mathbf{G}(s)$ with respect to $s$ evaluated at $s = \sigma$ is denoted by $\mathbf{G}^{(\ell)}(\sigma)$.

THEOREM 2.1. [1, 9] *Let $\sigma, \mu \in \mathbb{C}$ be such that $s\,\mathbf{E} - \mathbf{A}$ and $s\,\widetilde{\mathbf{E}} - \widetilde{\mathbf{A}}$ are both invertible for $s = \sigma, \mu$, and let $b \in \mathbb{C}^m$ and $c \in \mathbb{C}^p$ be fixed nontrivial vectors.*
1. *If*

$$
\left((\sigma\,\mathbf{E} - \mathbf{A})^{-1}\,\mathbf{E}\right)^{j-1}(\sigma\,\mathbf{E} - \mathbf{A})^{-1}\,\mathbf{B}b \in \mathrm{Im}(\mathbf{V}), \;\; j = 1, \ldots, N,
\tag{2.2}
$$

    *then $\mathbf{G}^{(\ell)}(\sigma)b = \widetilde{\mathbf{G}}^{(\ell)}(\sigma)b$ for $\ell = 0, 1, \ldots, N - 1$.*
2. *If*

$$
\left((\mu\,\mathbf{E} - \mathbf{A})^{-T}\,\mathbf{E}^T\right)^{j-1}(\mu\,\mathbf{E} - \mathbf{A})^{-T}\,\mathbf{C}^T c \in \mathrm{Im}(\mathbf{W}), \;\; j = 1, \ldots, M,
\tag{2.3}
$$

    *then $c^T \mathbf{G}^{(\ell)}(\mu) = c^T \widetilde{\mathbf{G}}^{(\ell)}(\mu)$ for $\ell = 0, 1, \ldots, M - 1$.*
3. *If both (2.2) and (2.3) hold, and if $\sigma = \mu$, then $c^T \mathbf{G}^{(\ell)}(\sigma)b = c^T \widetilde{\mathbf{G}}^{(\ell)}(\sigma)b$ for $\ell = 0, 1, \ldots, M + N + 1$.*

One can see that to solve the rational tangential interpolation problem via projection all one has to do is to construct the matrices $\mathbf{V}$ and $\mathbf{W}$ as in Theorem 2.1. The dominant cost is to solve sparse linear systems. We also note that in Theorem 2.1 the values that are interpolated are never explicitly computed. This is crucial since that computation is known to be poorly conditioned [8].

To illustrate the result of Theorem 2.1 for a special case of Hermite bi-tangential interpolation, we take the same right and left interpolation points $\{\sigma_i\}_{i=1}^r$, left tangential directions $\{c_i\}_{i=1}^r$, and right tangential directions $\{b_i\}_{i=1}^r$. Then for the projection matrices

$$
\mathbf{V} = \left[(\sigma_1\mathbf{E} - \mathbf{A})^{-1}\mathbf{B}b_1, \;\; \cdots, \;\; (\sigma_r\mathbf{E} - \mathbf{A})^{-1}\mathbf{B}b_r\right],
\tag{2.4}
$$

$$
\mathbf{W} = \left[(\sigma_1\mathbf{E} - \mathbf{A})^{-T}\mathbf{C}^T c_1, \;\; \cdots, \;\; (\sigma_r\mathbf{E} - \mathbf{A})^{-T}\mathbf{C}^T c_r\right],
\tag{2.5}
$$

the reduced-order model $\widetilde{\mathbf{G}}(s) = \widetilde{\mathbf{C}}(s\widetilde{\mathbf{E}} - \widetilde{\mathbf{A}})^{-1}\widetilde{\mathbf{B}} + \widetilde{\mathbf{D}}$ as in (1.3) satisfies

$$
\mathbf{G}(\sigma_i)b_i = \widetilde{\mathbf{G}}(\sigma_i)b_i, \quad c_i^T \mathbf{G}(\sigma_i) = c_i^T \widetilde{\mathbf{G}}(\sigma_i), \quad c_i^T \mathbf{G}'(\sigma_i)b_i = c_i^T \widetilde{\mathbf{G}}'(\sigma_i)b_i
\tag{2.6}
$$

for $i = 1, \cdots, r$, provided $\sigma_i\mathbf{E} - \mathbf{A}$ and $\sigma_i\widetilde{\mathbf{E}} - \widetilde{\mathbf{A}}$ are both nonsingular.

Note that Theorem 2.1 does not distinguish between the singular $\mathbf{E}$ case and the standard state space case with $\mathbf{E} = \mathbf{I}$. In other words, the interpolation conditions hold regardless as long as the matrices $\sigma_i\mathbf{E} - \mathbf{A}$ and $\sigma_i\widetilde{\mathbf{E}} - \widetilde{\mathbf{A}}$ are invertible. This is the precise reason why it is often assumed that extending interpolatory-based model reduction from $\mathbf{G}(s) = \mathbf{C}(s\mathbf{I}-\mathbf{A})^{-1}\mathbf{B}+\mathbf{D}$ to $\mathbf{G}(s) = \mathbf{C}(s\mathbf{E}-\mathbf{A})^{-1}\mathbf{B}+\mathbf{D}$ is as simple as replacing $\mathbf{I}$ by $\mathbf{E}$. However, as the following example shows, this is not the case.

EXAMPLE 2.1. Consider an RLC circuit modeled by an index-2 SISO descriptor system (1.1) of order $n = 765$ (see, e.g., [20] for a definition of index). We approximate this system with a model (1.2) of order $r = 20$ using Hermite interpolation. The carefully chosen interpolation points were taken as the mirror images of the dominant poles of $\mathbf{G}(s)$. Since these interpolation points are known to be good points for model reduction [11, 13], one would expect the interpolant to be a good approximation as well. However, the situation is indeed the opposite. Figure 2.1 shows the amplitude plots of the frequency responses $\mathbf{G}(\imath\omega)$ and $\widetilde{\mathbf{G}}(\imath\omega)$ (upper plot) and that of the error $\mathbf{G}(\imath\omega) - \widetilde{\mathbf{G}}(\imath\omega)$ (lower plot). One can see that the error $\mathbf{G}(\imath\omega) - \widetilde{\mathbf{G}}(\imath\omega)$ grows unbounded as the frequency $\omega$ increases, and, hence, the approximation is extremely poor with unbounded $\mathcal{H}_2$ and $\mathcal{H}_\infty$ error norms even though it satisfies Hermite interpolation at carefully selected effective interpolation points.
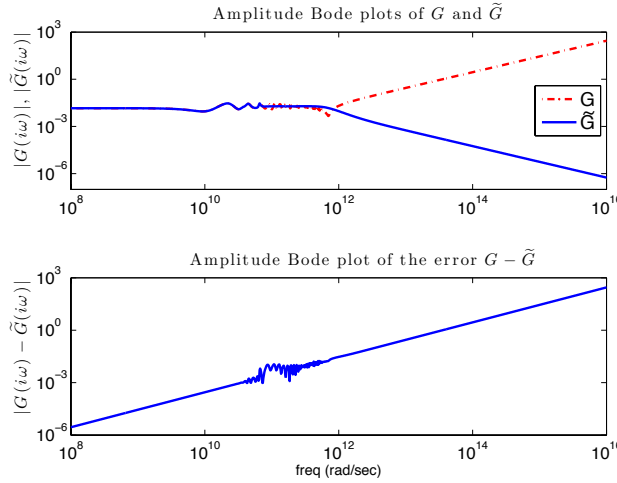


FIG. 2.1. *Example 2.1: amplitude plots of $\mathbf{G}(\imath\omega)$ and $\widetilde{\mathbf{G}}(\imath\omega)$ (upper); the absolute error $|\mathbf{G}(\imath\omega) - \widetilde{\mathbf{G}}(\imath\omega)|$ (lower).*

The reason is simple. Even though $\mathbf{E}$ is singular, $\widetilde{\mathbf{E}} = \mathbf{W}^T\mathbf{E}\mathbf{V}$ will generically be a nonsingular matrix assuming $r \leq \text{rank}(\mathbf{E})$. In this case, the transfer function $\widetilde{\mathbf{G}}(s)$ of the reduced-order model (1.2) is proper, i.e., $\lim_{s\to\infty} \widetilde{\mathbf{G}}(s) < \infty$, although $\mathbf{G}(s)$ might be improper. Hence, the special care needs to be taken in order to match the polynomial part of $\mathbf{G}(s)$. We note that the polynomial part of $\widetilde{\mathbf{G}}(s)$ has to match that of $\mathbf{G}(s)$ *exactly*. Otherwise, regardless of how good the interpolation points are, the error will always grow unbounded. For the very special descriptor systems with the proper transfer functions and only for interpolation around $s = \infty$, a solution is offered in [5]. For descriptor systems of index 1, where the polynomial part of $\mathbf{G}(s)$ is a constant matrix, a remedy is also suggested in [1] by an appropriate choice of $\widetilde{\mathbf{D}}$. However, the general case is remained unsolved. We will tackle precisely this problem,

where (1.1) is a descriptor system of higher index, its transfer function $\mathbf{G}(s)$ may have a higher order polynomial part and interpolation is at arbitrary points in the complex plane. Thereby, the spectral projectors onto the left and right deflating subspaces of the pencil $\lambda\mathbf{E} - \mathbf{A}$ corresponding to finite eigenvalues will play a vital role. Moreover, we will show how to choose interpolation points and tangential directions optimally for interpolatory model reduction of descriptor systems.

**3. Interpolatory projection methods for descriptor systems.** As stated above, in order to have bounded $\mathcal{H}_\infty$ and $\mathcal{H}_2$ errors, the polynomial part of $\widetilde{\mathbf{G}}(s)$ has to match the polynomial part of $\mathbf{G}(s)$ exactly. Let $\mathbf{G}(s)$ be additively decomposed as

$$\mathbf{G}(s) = \mathbf{G}_{\mathrm{sp}}(s) + \mathbf{P}(s), \tag{3.1}$$

where $\mathbf{G}_{\mathrm{sp}}(s)$ and $\mathbf{P}(s)$ denote, respectively, the strictly proper part and the polynomial part of $\mathbf{G}(s)$. We enforce the reduced-order model $\widetilde{\mathbf{G}}(s)$ to have the decomposition

$$\widetilde{\mathbf{G}}(s) = \widetilde{\mathbf{G}}_{\mathrm{sp}}(s) + \widetilde{\mathbf{P}}(s) \tag{3.2}$$

with $\widetilde{\mathbf{P}}(s) = \mathbf{P}(s)$. This implies that the error transfer function does not contain a polynomial part, i.e.,

$$\mathbf{G}_{\mathrm{err}}(s) = \mathbf{G}(s) - \widetilde{\mathbf{G}}(s) = \mathbf{G}_{\mathrm{sp}}(s) - \widetilde{\mathbf{G}}_{\mathrm{sp}}(s)$$

is strictly proper meaning $\lim_{s\to\infty} \mathbf{G}_{\mathrm{err}}(s) = 0$. Hence, by making $\widetilde{\mathbf{G}}_{\mathrm{sp}}(s)$ to interpolate $\mathbf{G}_{\mathrm{sp}}(s)$, we will be able to enforce that $\widetilde{\mathbf{G}}(s)$ interpolates $\mathbf{G}(s)$. This will lead to the following construction of $\widetilde{\mathbf{G}}(s)$. Given $\mathbf{G}(s)$, we create $\mathbf{W}$ and $\mathbf{V}$ satisfying new subspace conditions such that the reduced-order model $\widetilde{\mathbf{G}}(s)$ obtained by projection as in (1.3) will not only satisfy the interpolation conditions but also match the polynomial part of $\mathbf{G}(s)$.

THEOREM 3.1. *Given a full-order model* $\mathbf{G}(s) = \mathbf{C}(s\mathbf{E} - \mathbf{A})^{-1}\mathbf{B} + \mathbf{D}$, *define* $\mathbf{P}_l$ *and* $\mathbf{P}_r$ *to be the spectral projectors onto the left and right deflating subspaces of the pencil* $\lambda\mathbf{E} - \mathbf{A}$ *corresponding to the finite eigenvalues. Let the columns of* $\mathbf{W}_\infty$ *and* $\mathbf{V}_\infty$ *span the left and right deflating subspaces of* $\lambda\mathbf{E} - \mathbf{A}$ *corresponding to the eigenvalue at infinity. Let* $\sigma$, $\mu \in \mathbb{C}$ *be interpolation points such that* $s\mathbf{E} - \mathbf{A}$ *and* $s\widetilde{\mathbf{E}} - \widetilde{\mathbf{A}}$ *are nonsingular for* $s = \sigma, \mu$, *and let* $\mathsf{b} \in \mathbb{C}^m$ *and* $\mathsf{c} \in \mathbb{C}^p$. *Define* $\mathbf{V}_f$ *and* $\mathbf{W}_f$ *such that*

$$\mathrm{Im}(\mathbf{V}_f) = \mathrm{span}\left\{ \left((\sigma\mathbf{E} - \mathbf{A})^{-1}\mathbf{E}\right)^{j-1} (\sigma\mathbf{E} - \mathbf{A})^{-1}\mathbf{P}_l\mathbf{B}\mathsf{b}, \ \ j = 1, ..., N \right\}, \tag{3.3}$$

$$\mathrm{Im}(\mathbf{W}_f) = \mathrm{span}\left\{ \left((\mu\mathbf{E} - \mathbf{A})^{-T}\mathbf{E}^T\right)^{j-1} (\mu\mathbf{E} - \mathbf{A})^{-T}\mathbf{P}_r^T\mathbf{C}^T\mathsf{c}, \ \ j = 1, ..., M \right\}. \tag{3.4}$$

*Then with the choice of* $\mathbf{W} = [\,\mathbf{W}_f, \ \mathbf{W}_\infty\,]$ *and* $\mathbf{V} = [\,\mathbf{V}_f, \ \mathbf{V}_\infty\,]$, *the reduced-order model* $\widetilde{\mathbf{G}}(s) = \widetilde{\mathbf{C}}(s\widetilde{\mathbf{E}} - \widetilde{\mathbf{A}})^{-1}\widetilde{\mathbf{B}} + \widetilde{\mathbf{D}}$ *obtained via projection as in (1.3) satisfies*
   1. *$\widetilde{\mathbf{P}}(s) = \mathbf{P}(s)$,*
   2. *$\mathbf{G}^{(\ell)}(\sigma)\mathsf{b} = \widetilde{\mathbf{G}}^{(\ell)}(\sigma)\mathsf{b}$ for $\ell = 0, 1, \ldots, N-1$,*
   3. *$\mathsf{c}^T\mathbf{G}^{(\ell)}(\mu) = \mathsf{c}^T\widetilde{\mathbf{G}}^{(\ell)}(\mu)$ for $\ell = 0, 1, \ldots, M-1$.*
*If $\sigma = \mu$, we have, additionally, $\mathsf{c}^T\mathbf{G}^{(\ell)}(\sigma)\mathsf{b} = \mathsf{c}^T\widetilde{\mathbf{G}}^{(\ell)}(\sigma)\mathsf{b}$ for $\ell = 0, \ldots, M+N+1$.*
   *Proof.* Let the pencil $\lambda\mathbf{E} - \mathbf{A}$ be transformed into the Weierstrass canonical form

$$\mathbf{E} = \mathbf{S} \begin{bmatrix} \mathbf{I}_{n_f} & \mathbf{0} \\ \mathbf{0} & \mathbf{N} \end{bmatrix} \mathbf{T}^{-1}, \qquad \mathbf{A} = \mathbf{S} \begin{bmatrix} \mathbf{J} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n_\infty} \end{bmatrix} \mathbf{T}^{-1}, \tag{3.5}$$

where $\mathbf{S}$ and $\mathbf{T}$ are nonsingular and $\mathbf{N}$ is nilpotent. Then the projectors $\mathbf{P}_l$ and $\mathbf{P}_r$ can be represented as

$$\mathbf{P}_l = \mathbf{S} \begin{bmatrix} \mathbf{I}_{n_f} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{S}^{-1}, \qquad \mathbf{P}_r = \mathbf{T} \begin{bmatrix} \mathbf{I}_{n_f} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{T}^{-1}. \tag{3.6}$$

Let $\mathbf{T} = [\mathbf{T}_1, \mathbf{T}_2]$ and $\mathbf{S}^{-1} = [\mathbf{S}_1, \mathbf{S}_2]^T$ be partitioned according to $\mathbf{E}$ and $\mathbf{A}$ in (3.5). Then the matrices $\mathbf{W}_\infty$ and $\mathbf{V}_\infty$ take the form

$$\mathbf{W}_\infty = \mathbf{S}_2 \mathbf{R}_S = (\mathbf{I} - \mathbf{P}_l^T)\mathbf{W}_\infty, \qquad \mathbf{V}_\infty = \mathbf{T}_2 \mathbf{R}_T = (\mathbf{I} - \mathbf{P}_r)\mathbf{V}_\infty$$

with nonsingular $\mathbf{R}_S$ and $\mathbf{R}_T$. Furthermore, the strictly proper and polynomial parts of $\mathbf{G}(s)$ in (3.1) are given by

$$\mathbf{G}_{\mathrm{sp}}(s) = \mathbf{C}\mathbf{T}_1(s\mathbf{I}_{n_f} - \mathbf{J})^{-1}\mathbf{S}_1^T\mathbf{B},$$

$$\text{and} \quad \mathbf{P}(s) = \mathbf{C}\mathbf{T}_2(s\mathbf{N} - \mathbf{I}_{n_\infty})^{-1}\mathbf{S}_2^T\mathbf{B} + \mathbf{D},$$

respectively. It follows from (3.5) and (3.6) that

$$\mathbf{E}\mathbf{P}_r = \mathbf{P}_l\mathbf{E}, \qquad \mathbf{A}\mathbf{P}_r = \mathbf{P}_l\mathbf{A},$$
$$(s\mathbf{E} - \mathbf{A})^{-1}\mathbf{P}_l = \mathbf{P}_r(s\mathbf{E} - \mathbf{A})^{-1},$$

and, hence,

$$\mathbf{W}_f = \mathbf{P}_l^T\mathbf{W}_f \qquad \text{and} \qquad \mathbf{V}_f = \mathbf{P}_r\mathbf{V}_f. \tag{3.7}$$

Then the system matrices of the reduced-order model have the form

$$\widetilde{\mathbf{E}} = \mathbf{W}^T\mathbf{E}\mathbf{V} = \begin{bmatrix} \mathbf{W}_f^T\mathbf{E}\mathbf{V}_f & \mathbf{W}_f^T\mathbf{E}\mathbf{V}_\infty \\ \mathbf{W}_\infty^T\mathbf{E}\mathbf{V}_f & \mathbf{W}_\infty^T\mathbf{E}\mathbf{V}_\infty \end{bmatrix} = \begin{bmatrix} \mathbf{W}_f^T\mathbf{E}\mathbf{V}_f & \mathbf{0} \\ \mathbf{0} & \mathbf{W}_\infty^T\mathbf{E}\mathbf{V}_\infty \end{bmatrix},$$

$$\widetilde{\mathbf{A}} = \mathbf{W}^T\mathbf{A}\mathbf{V} = \begin{bmatrix} \mathbf{W}_f^T\mathbf{A}\mathbf{V}_f & \mathbf{W}_f^T\mathbf{A}\mathbf{V}_\infty \\ \mathbf{W}_\infty^T\mathbf{A}\mathbf{V}_f & \mathbf{W}_\infty^T\mathbf{A}\mathbf{V}_\infty \end{bmatrix} = \begin{bmatrix} \mathbf{W}_f^T\mathbf{A}\mathbf{V}_f & \mathbf{0} \\ \mathbf{0} & \mathbf{W}_\infty^T\mathbf{A}\mathbf{V}_\infty \end{bmatrix},$$

$$\widetilde{\mathbf{B}} = \mathbf{W}^T\mathbf{B} = \begin{bmatrix} \mathbf{W}_f^T\mathbf{B} \\ \mathbf{W}_\infty^T\mathbf{B} \end{bmatrix}, \qquad \widetilde{\mathbf{C}} = \mathbf{C}\mathbf{V} = [\mathbf{C}\mathbf{V}_f, \mathbf{C}\mathbf{V}_\infty], \qquad \widetilde{\mathbf{D}} = \mathbf{D}.$$

Thus, the strictly proper and polynomial parts of $\widetilde{\mathbf{G}}(s)$ are given by

$$\widetilde{\mathbf{G}}_{\mathrm{sp}}(s) = \mathbf{C}\mathbf{V}_f(s\mathbf{W}_f^T\mathbf{E}\mathbf{V}_f - \mathbf{W}_f^T\mathbf{A}\mathbf{V}_f)^{-1}\mathbf{W}_f^T\mathbf{B},$$

$$\widetilde{\mathbf{P}}(s) = \mathbf{C}\mathbf{V}_\infty(s\mathbf{W}_\infty^T\mathbf{E}\mathbf{V}_\infty - \mathbf{W}_\infty^T\mathbf{A}\mathbf{V}_\infty)^{-1}\mathbf{W}_\infty^T\mathbf{B} + \mathbf{D}$$

$$= \mathbf{C}\mathbf{T}_2(s\mathbf{I} - \mathbf{J})^{-1}\mathbf{S}_2^T\mathbf{B} + \mathbf{D} = \mathbf{P}(s).$$

One can see that the polynomial parts of $\mathbf{G}(s)$ and $\widetilde{\mathbf{G}}(s)$ coincide, and the proof of the interpolation result reduces to proving the interpolation conditions for the strictly proper parts of $\mathbf{G}(s)$ and $\widetilde{\mathbf{G}}(s)$. To prove this, we first note that (3.5) and (3.6) imply that

$$\mathbf{C}\mathbf{P}_r(\sigma\mathbf{E} - \mathbf{A})^{-1}\mathbf{P}_l\mathbf{B} = \mathbf{C}\mathbf{T}\begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}\begin{bmatrix} \sigma\mathbf{I} - \mathbf{J} & \mathbf{0} \\ \mathbf{0} & \sigma\mathbf{N} - \mathbf{I} \end{bmatrix}^{-1}\begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}\mathbf{S}^{-1}\mathbf{B}$$

$$= \mathbf{C}\mathbf{T}_1(\sigma\mathbf{I} - \mathbf{J})^{-1}\mathbf{S}_1^T\mathbf{B} = \mathbf{G}_{\mathrm{sp}}(\sigma).$$

Furthermore, it follows from the relations (3.7) that

$$\mathbf{CP}_r\mathbf{V}_f = \mathbf{CV}_f, \qquad \mathbf{W}_f^T\mathbf{P}_l\mathbf{B} = \mathbf{W}_f^T\mathbf{B}.$$

Due to the definitions of $\mathbf{V}_f$ and $\mathbf{W}_f$ in (3.3) and (3.4), respectively, Theorem 2.1 gives

$$\mathbf{G}_{\mathrm{sp}}(\sigma)\mathsf{b} = \mathbf{CP}_r\mathbf{V}_f(\sigma\mathbf{W}_f^T\mathbf{EV}_f - \mathbf{W}_f^T\mathbf{AV}_f)^{-1}\mathbf{W}_f^T\mathbf{P}_l\mathbf{Bb} = \widetilde{\mathbf{G}}_{\mathrm{sp}}(\sigma)\mathsf{b},$$
$$\mathsf{c}^T\mathbf{G}_{\mathrm{sp}}(\mu) = \mathsf{c}^T\mathbf{CP}_r\mathbf{V}_f(\mu\mathbf{W}_f^T\mathbf{EV}_f - \mathbf{W}_f^T\mathbf{AV}_f)^{-1}\mathbf{W}_f^T\mathbf{P}_l\mathbf{B} = \mathsf{c}^T\widetilde{\mathbf{G}}_{\mathrm{sp}}(\mu).$$

Since both parts 1 and 2 of Theorem 2.1 hold, we have $\mathsf{c}^T\widetilde{\mathbf{G}}'_{\mathrm{sp}}(\sigma)\mathsf{b} = \mathsf{c}^T\mathbf{G}'_{\mathrm{sp}}(\sigma)\mathsf{b}$ for $\sigma = \mu$. The other interpolatory relations for the derivatives of the transfer function can be proved analogously. □

Next, we illustrate that even though Theorem 3.1 has a very similar structure to that of Theorem 2.1, the saddle difference between these two results makes a big difference in the resulting reduced-order model. Towards this goal, we revisit Example 2.1. We reduce the same full-order model using the same interpolation points, but imposing the subspace conditions of Theorem 3.1, instead. Figure 3.1 depicts the resulting amplitude plots of $\mathbf{G}(\imath\omega)$ and $\widetilde{\mathbf{G}}(\imath\omega)$ (upper plot) and that of the error $\mathbf{G}(\imath\omega) - \widetilde{\mathbf{G}}(\imath\omega)$ (lower plot) when the new subspace conditions of Theorem 3.1 are used. Unlike the case in Example 2.1, where the error $\mathbf{G}(\imath\omega) - \widetilde{\mathbf{G}}(\imath\omega)$ grew unbounded, for the new reduced-order model, the maximum error is below $10^{-2}$ and the error decays to zero as $\omega$ approaches $\infty$, since the polynomial part is captured exactly.



FIG. 3.1. *Amplitude plots of* $\mathbf{G}(\imath\omega)$ *and* $\widetilde{\mathbf{G}}(\imath\omega)$ *(upper); the absolute error* $|\mathbf{G}(\imath\omega) - \widetilde{\mathbf{G}}(\imath\omega)|$ *(lower).*

In some applications, the deflating subspaces of $\lambda\mathbf{E} - \mathbf{A}$ corresponding to the eigenvalues at infinity may have large dimension $n_\infty$. However, the order of the system can still be reduced if it contains states that are uncontrollable and unobservable at infinity. Such states can be removed from the system without changing its transfer

function and, hence, preserving the interpolation conditions as in Theorem 3.1. In this case the projection matrices $\mathbf{W}_\infty$ and $\mathbf{V}_\infty$ can be determined as proposed in [26] by solving the projected discrete-time Lyapunov equations

$$\mathbf{A}\mathbf{X}\mathbf{A}^T - \mathbf{E}\mathbf{X}\mathbf{E}^T = -(\mathbf{I} - \mathbf{P}_l)\mathbf{B}\mathbf{B}^T(\mathbf{I} - \mathbf{P}_l)^T, \quad \mathbf{X} = (\mathbf{I} - \mathbf{P}_r)\mathbf{X}(\mathbf{I} - \mathbf{P}_r)^T, \quad (3.8)$$

$$\mathbf{A}^T\mathbf{Y}\mathbf{A} - \mathbf{E}^T\mathbf{Y}\mathbf{E} = -(\mathbf{I} - \mathbf{P}_r)^T\mathbf{C}^T\mathbf{C}(\mathbf{I} - \mathbf{P}_r), \quad \mathbf{Y} = (\mathbf{I} - \mathbf{P}_l)^T\mathbf{Y}(\mathbf{I} - \mathbf{P}_l). \quad (3.9)$$

Let $\mathbf{X}_C$ and $\mathbf{Y}_C$ be the Cholesky factors of $\mathbf{X} = \mathbf{X}_C\mathbf{X}_C^T$ and $\mathbf{Y} = \mathbf{Y}_C\mathbf{Y}_C^T$, respectively, and let $\mathbf{Y}_C^T\mathbf{A}\mathbf{X}_C = [\mathbf{U}_1, \mathbf{U}_0]\mathrm{diag}(\mathbf{\Sigma}, \mathbf{0})[\mathbf{V}_1, \mathbf{V}_0]^T$ be singular value decomposition, where $[\mathbf{U}_1, \mathbf{U}_0]$ and $[\mathbf{V}_1, \mathbf{V}_0]$ are orthogonal and $\mathbf{\Sigma}$ is nonsingular. Then the projection matrices $\mathbf{W}_\infty$ and $\mathbf{V}_\infty$ can be taken as $\mathbf{W}_\infty = \mathbf{Y}_C\mathbf{U}_1$ and $\mathbf{V}_\infty = \mathbf{X}_C\mathbf{V}_1$. Note that the Cholesky factors $\mathbf{X}_C$ and $\mathbf{Y}_C$ can be computed directly using the generalized Smith method [27]. In this method, it is required to solve $\nu - 1$ linear systems only, where $\nu$ is the index of the pencil $\lambda\mathbf{E} - \mathbf{A}$ or, equivalently, the nilpotence index of $\mathbf{N}$ in (3.5). The computation of the projectors $\mathbf{P}_l$ and $\mathbf{P}_r$ is, in general, a difficult problem. However, for some structured problems arising in circuit simulation, multibody systems and computational fluid dynamics, these projectors can be constructed in explicit form that significantly reduces the computational complexity of the method; see [27] for details.

**4. Interpolatory optimal $\mathcal{H}_2$ model reduction for descriptor systems.** The choice of interpolation points and tangential directions is the central issue in interpolatory model reduction. This choice determines whether the reduced-order model is high fidelity or not. Until recently, selection of interpolation points was largely *ad hoc* and required several model reduction attempts to arrive at a reasonable approximation. However, Gugercin *et al.* [15] introduced an interpolatory model reduction method for generating a reduced model $\widetilde{\mathbf{G}}$ of order $r$ which is an optimal $\mathcal{H}_2$ approximation to the original system $\mathbf{G}$ in the sense that it minimizes $\mathcal{H}_2$-norm error, i.e.,

$$\|\mathbf{G} - \widetilde{\mathbf{G}}\|_{\mathcal{H}_2} = \min_{\dim(\widetilde{\mathbf{G}}_r)=r} \|\mathbf{G} - \widetilde{\mathbf{G}}_r\|_{\mathcal{H}_2}, \quad (4.1)$$

where

$$\|\mathbf{G}\|_{\mathcal{H}_2} := \left(\frac{1}{2\pi}\int_{-\infty}^{+\infty}\|\mathbf{G}(\imath\omega)\|_{\mathrm{F}}^2\,d\omega\right)^{1/2} \quad (4.2)$$

and $\|\cdot\|_{\mathrm{F}}$ denotes the Frobenius matrix norm. Since this is a non-convex optimization problem, the computation of a global minimizer is a very difficult task. Hence, instead, one tries to find high-fidelity reduced models that satisfy first-order necessary optimality conditions. There exist, in general, two approaches for solving this problem. These are Lyapunov-based optimal $\mathcal{H}_2$ methods presented in [16, 18, 25, 29, 30, 33] and interpolation-based optimal $\mathcal{H}_2$ methods considered in [3, 4, 6, 12, 14, 15, 19, 23, 28]. While the Lyapunov-based approaches require solving a series of Lyapunov equations, which becomes costly and sometimes intractable in large-scale settings, the interpolatory approaches only require solving a series of sparse linear systems and have proved to be numerically very effective. Moreover, as shown in [15], both frameworks are theoretically equivalent that further motivates the usage of interpolatory model reduction techniques for the optimal $\mathcal{H}_2$ approximation.

For SISO systems, interpolation-based $\mathcal{H}_2$ optimality conditions were originally developed by Meier and Luenberger [23]. Then, based on these conditions, an effective

algorithm for interpolatory optimal $\mathcal{H}_2$ approximation, called the *Iterative Rational Krylov Algorithm* (IRKA), was introduced in [12, 14]. This algorithm has also been recently extended to MIMO systems using the tangential interpolation framework, see [6, 15, 28] for more details.

The model reduction methods mentioned above, however, only deals with the system $\mathbf{G}(s) = \mathbf{C}(s\mathbf{E} - \mathbf{A})^{-1}\mathbf{B} + \mathbf{D}$ with a nonsingular matrix $\mathbf{E}$. In this section, we will extend IRKA to descriptor systems. First, we establish the interpolatory $\mathcal{H}_2$ optimality conditions in the new setting.

THEOREM 4.1. *Let* $\mathbf{G}(s) = \mathbf{G}_{\text{sp}}(s) + \mathbf{P}(s)$ *be decomposed into the strictly proper and polynomial parts, and let* $\widetilde{\mathbf{G}}(s) = \widetilde{\mathbf{G}}_{\text{sp}}(s) + \widetilde{\mathbf{P}}(s)$ *have an* $r^{\text{th}}$*-order strictly proper part* $\widetilde{\mathbf{G}}_{\text{sp}}(s) = \widetilde{\mathbf{C}}_{\text{sp}}(s\widetilde{\mathbf{E}}_{\text{sp}} - \widetilde{\mathbf{A}}_{\text{sp}})^{-1}\widetilde{\mathbf{B}}_{\text{sp}}$.

1. *If* $\widetilde{\mathbf{G}}(s)$ *minimizes the* $\mathcal{H}_2$*-error* $\|\mathbf{G} - \widetilde{\mathbf{G}}\|_{\mathcal{H}_2}$ *over all reduced-order models with an* $r^{\text{th}}$*-order strictly proper part, then* $\widetilde{\mathbf{P}}(s) = \mathbf{P}(s)$ *and* $\widetilde{\mathbf{G}}_{\text{sp}}(s)$ *minimizes the* $\mathcal{H}_2$*-error* $\|\mathbf{G}_{\text{sp}} - \widetilde{\mathbf{G}}_{\text{sp}}\|_{\mathcal{H}_2}$.

2. *Suppose that the pencil* $s\widetilde{\mathbf{E}}_{\text{sp}} - \widetilde{\mathbf{A}}_{\text{sp}}$ *has distinct eigenvalues* $\{\widetilde{\lambda}_i\}_{i=1}^r$. *Let* $\mathbf{y}_i$ *and* $\mathbf{z}_i$ *denote the left and right eigenvectors associated with* $\widetilde{\lambda}_i$ *so that* $\widetilde{\mathbf{A}}_{\text{sp}}\mathbf{z}_i = \widetilde{\lambda}_i\widetilde{\mathbf{E}}_{\text{sp}}\mathbf{z}_i$, $\mathbf{y}_i^*\widetilde{\mathbf{A}}_{\text{sp}} = \widetilde{\lambda}_i\mathbf{y}_i^*\widetilde{\mathbf{E}}_{\text{sp}}$, *and* $\mathbf{y}_i^*\widetilde{\mathbf{E}}_{\text{sp}}\mathbf{z}_j = \delta_{ij}$. *Then for* $\mathsf{c}_i = \widetilde{\mathbf{C}}_{\text{sp}}\mathbf{z}_i$ *and* $\mathsf{b}_i^T = \mathbf{y}_i^*\widetilde{\mathbf{B}}_{\text{sp}}$, *we have*

$$\mathbf{G}(-\widetilde{\lambda}_i)\mathsf{b}_i = \widetilde{\mathbf{G}}(-\widetilde{\lambda}_i)\mathsf{b}_i, \qquad \mathsf{c}_i^T\mathbf{G}(-\widetilde{\lambda}_i) = \mathsf{c}_i^T\widetilde{\mathbf{G}}(-\widetilde{\lambda}_i),$$
$$\text{and} \quad \mathsf{c}_i^T\mathbf{G}'(-\widetilde{\lambda}_i)\mathsf{b}_i = \mathsf{c}_i^T\widetilde{\mathbf{G}}'(-\widetilde{\lambda}_i)\mathsf{b}_i \quad \text{for } i = 1, \cdots, r. \tag{4.3}$$

*Proof.* 1. The polynomial part of $\mathbf{G}(s)$ and $\widetilde{\mathbf{G}}(s)$ coincide, since, otherwise, the $\mathcal{H}_2$-norm of the error $\mathbf{G}(s) - \widetilde{\mathbf{G}}(s)$ would be unbounded. Then it readily follows that $\widetilde{\mathbf{G}}_{\text{sp}}(s)$ minimizes $\|\mathbf{G}_{\text{sp}}(s) - \widetilde{\mathbf{G}}_{\text{sp}}(s)\|_{\mathcal{H}_2}$ since $\mathbf{G}(s) - \widetilde{\mathbf{G}}(s) = \mathbf{G}_{\text{sp}}(s) - \widetilde{\mathbf{G}}_{\text{sp}}(s)$.

2. Since $\widetilde{\mathbf{P}}(s) = \mathbf{P}(s)$, the $\mathcal{H}_2$ optimal model reduction problem for $\mathbf{G}(s)$ now reduces to the $\mathcal{H}_2$ optimal problem for the strictly proper transfer function $\mathbf{G}_{\text{sp}}(s)$. Hence, the optimal $\mathcal{H}_2$ conditions of [15] require that $\widetilde{\mathbf{G}}_{\text{sp}}(s)$ needs to be a bi-tangential Hermite interpolant to $\mathbf{G}_{\text{sp}}(s)$ with $\{-\widetilde{\lambda}_i\}_{i=1}^r$ being the interpolation points, and $\{\mathsf{c}_i\}_{i=1}^r$ and $\{\mathsf{b}_i\}_{i=1}^r$ being the corresponding left and right tangential directions, respectively. Thus, the interpolation conditions (4.3) hold since $\widetilde{\mathbf{P}}(s) = \mathbf{P}(s)$. $\square$

Unfortunately, the $\mathcal{H}_2$ optimal interpolation points and associated tangent directions are not known *a priori*, since they depend on the reduced-order model to be computed. To overcome this difficulty, an iterative algorithm IRKA was developed [12, 14] which is based on successive substitution. In IRKA, the interpolation points are corrected iteratively by the choosing mirror images of poles of the current reduced-order model as the next interpolation points. The tangential directions are corrected in a similar way; see [1, 15] for details.

The situation in the case of descriptor systems is similar, where the optimal interpolation points and the corresponding tangential directions depend on the strictly proper part of the reduced-order model to be computed. Moreover, we need to make sure that the final reduced-model has the same polynomial part as the original one. Hence, we will modify IRKA to meet these challenges. In particular, we will correct not the poles and the tangential directions of the intermediate reduced-order model at the successive iteration step but that of the strictly proper part of the intermediate reduced-order model. As in the case of Theorem 3.1, the spectral projectors $\mathbf{P}_l$ and $\mathbf{P}_r$ will be used to construct the required interpolatory subspaces. A sketch of the resulting model reduction method is given in Algorithm 4.1.

---

ALGORITHM 4.1. **Interpolatory $\mathcal{H}_2$ optimal model reduction method for descriptor systems**

1) *Make an initial selection of the interpolation points $\{\sigma_i\}_{i=1}^r$ and the tangential directions $\{\mathsf{b}_i\}_{i=1}^r$ and $\{\mathsf{c}_i\}_{i=1}^r$.*

2) $\mathbf{V}_f = \left[\, (\sigma_1 \mathbf{E} - \mathbf{A})^{-1} \mathbf{P}_l \mathbf{B} \mathsf{b}_1, \ \ldots, \ (\sigma_r \mathbf{E} - \mathbf{A})^{-1} \mathbf{P}_l \mathbf{B} \mathsf{b}_r \,\right]$,

   $\mathbf{W}_f = \left[\, (\sigma_1 \mathbf{E} - \mathbf{A})^{-T} \mathbf{P}_r^T \mathbf{C}^T \mathsf{c}_1, \ \ldots, \ (\sigma_r \mathbf{E} - \mathbf{A})^{-T} \mathbf{P}_r^T \mathbf{C}^T \mathsf{c}_r \,\right]$.

3) *while (not converged)*

   a) $\widetilde{\mathbf{A}}_{\mathrm{sp}} = \mathbf{W}_f^T \mathbf{A} \mathbf{V}_f$, $\widetilde{\mathbf{E}}_{\mathrm{sp}} = \mathbf{W}_f^T \mathbf{E} \mathbf{V}_f$, $\widetilde{\mathbf{B}}_{\mathrm{sp}} = \mathbf{W}_f^T \mathbf{B}$, *and* $\widetilde{\mathbf{C}}_{\mathrm{sp}} = \mathbf{C} \mathbf{V}_f$.

   b) *Compute* $\widetilde{\mathbf{A}}_{\mathrm{sp}} \mathbf{z}_i = \widetilde{\lambda}_i \widetilde{\mathbf{E}}_{\mathrm{sp}} \mathbf{z}_i$ *and* $\mathbf{y}_i^* \widetilde{\mathbf{A}}_{\mathrm{sp}} = \widetilde{\lambda}_i \mathbf{y}_i^* \widetilde{\mathbf{E}}_{\mathrm{sp}}$ *with* $\mathbf{y}_i^* \widetilde{\mathbf{E}}_{\mathrm{sp}} \mathbf{z}_j = \delta_{ij}$,

   *where $\mathbf{y}_i$ and $\mathbf{z}_i$ are left and right eigenvectors associated with $\widetilde{\lambda}_i$.*

   c) $\sigma_i \leftarrow -\widetilde{\lambda}_i$, $\mathsf{b}_i^T \leftarrow \mathbf{y}_i^* \widetilde{\mathbf{B}}_{\mathrm{sp}}$ *and* $\mathsf{c}_i \leftarrow \widetilde{\mathbf{C}}_{\mathrm{sp}} \mathbf{z}_i$ *for* $i = 1, \ldots, r$.

   d) $\mathbf{V}_f = \left[\, (\sigma_1 \mathbf{E} - \mathbf{A})^{-1} \mathbf{P}_l \mathbf{B} \mathsf{b}_1, \ \ldots, \ (\sigma_r \mathbf{E} - \mathbf{A})^{-1} \mathbf{P}_l \mathbf{B} \mathsf{b}_r \,\right]$,

   $\mathbf{W}_f = \left[\, (\sigma_1 \mathbf{E} - \mathbf{A})^{-T} \mathbf{P}_r^T \mathbf{C}^T \mathsf{c}_1, \ \ldots, \ (\sigma_r \mathbf{E} - \mathbf{A})^{-T} \mathbf{P}_r^T \mathbf{C}^T \mathsf{c}_r \,\right]$.

   *end while*

4) *Compute $\mathbf{W}_\infty$ and $\mathbf{V}_\infty$ such that* $\mathrm{Im}(\mathbf{W}_\infty) = \mathrm{Im}(\mathbf{I} - \mathbf{P}_l^T)$ *and* $\mathrm{Im}(\mathbf{V}_\infty) = \mathrm{Im}(\mathbf{I} - \mathbf{P}_r)$.

5) *Set* $\mathbf{V} = [\, \mathbf{V}_f, \ \mathbf{V}_\infty \,]$ *and* $\mathbf{W} = [\, \mathbf{W}_f, \ \mathbf{W}_\infty \,]$.

6) $\widetilde{\mathbf{E}} = \mathbf{W}^T \mathbf{E} \mathbf{V}$, $\widetilde{\mathbf{A}} = \mathbf{W}^T \mathbf{A} \mathbf{V}$, $\widetilde{\mathbf{B}} = \mathbf{W}^T \mathbf{B}$, $\widetilde{\mathbf{C}} = \mathbf{C} \mathbf{V}$, $\widetilde{\mathbf{D}} = \mathbf{D}$.

---

Note that until Step 4 of Algorithm 4.1, the polynomial part is not included since the interpolation parameters result from the strictly proper part $\widetilde{\mathbf{G}}_{\mathrm{sp}}(s)$. In a sense, Step 3 runs the optimal $\mathcal{H}_2$ iteration on $\mathbf{G}_{\mathrm{sp}}(s)$. Hence, at the end of Step 3, we construct an optimal $\mathcal{H}_2$ interpolant to $\mathbf{G}_{\mathrm{sp}}(s)$. However, in Step 5, we append the interpolatory subspaces with $\mathbf{V}_\infty$ and $\mathbf{W}_\infty$ (which can be computed as described at the end of Section 3) so that the final reduced-order model in Step 6 has the same polynomial part as $\mathbf{G}(s)$, and, consequently, the final reduced-order model $\widetilde{\mathbf{G}}(s)$ satisfies the optimality conditions of Theorem 4.1. One can see this from Step 3c: upon convergence, the interpolation points are the mirror images of the poles of $\widetilde{\mathbf{G}}_{\mathrm{sp}}(s)$ and the tangential directions are the residue directions from $\widetilde{\mathbf{G}}_{\mathrm{sp}}(s)$ as the optimality conditions require. Since Algorithm 4.1 uses the projected quantities $\mathbf{P}_l \mathbf{B}$ and $\mathbf{C} \mathbf{P}_r$, theoretically iterating on a strictly proper dynamical system, the convergence behavior of this algorithm will follow the same pattern of IRKA which has been observed to converge rapidly in numerous numerical applications.

Summarizing, we have shown so far how to reduce descriptor systems such that the transfer function of the reduced descriptor systems is a tangential interpolant to the original one and matches the polynomial part preventing unbounded $\mathcal{H}_\infty$ and $\mathcal{H}_2$ error norms. However, this model reduction approach involves the explicit computation of the spectral projectors or the corresponding deflating subspaces, which could be numerically infeasible for general large-scale problems. In the next two sections, we will show that for certain important classes of descriptor systems, the same can be achieved without explicitly forming the spectral projectors.

**5. Semi-explicit descriptor systems of index 1.** We consider the following semi-explicit descriptor system

$$
\begin{aligned}
\mathbf{E}_{11} \dot{\mathbf{x}}_1(t) + \mathbf{E}_{12} \dot{\mathbf{x}}_2(t) &= \mathbf{A}_{11} \mathbf{x}_1(t) + \mathbf{A}_{12} \mathbf{x}_2(t) + \mathbf{B}_1 \mathbf{u}(t), \\
\mathbf{0} &= \mathbf{A}_{21} \mathbf{x}_1(t) + \mathbf{A}_{22} \mathbf{x}_2(t) + \mathbf{B}_2 \mathbf{u}(t), \\
\mathbf{y}(t) &= \mathbf{C}_1 \mathbf{x}_1(t) + \mathbf{C}_2 \mathbf{x}_2(t) + \mathbf{D} \mathbf{u}(t),
\end{aligned}
\tag{5.1}
$$

where the state is $\mathbf{x}(t) = [\,\mathbf{x}_1^T(t),\ \mathbf{x}_2^T(t)\,]^T \in \mathbb{R}^n$ with $\mathbf{x}_1(t) \in \mathbb{R}^{n_1}$, $\mathbf{x}_2(t) \in \mathbb{R}^{n_2}$ and $n_1 + n_2 = n$, the input is $\mathbf{u}(t) \in \mathbb{R}^m$, the output is $\mathbf{y}(t) \in \mathbb{R}^p$, and $\mathbf{E}_{11}, \mathbf{A}_{11} \in \mathbb{R}^{n_1 \times n_1}$, $\mathbf{E}_{12}, \mathbf{A}_{12} \in \mathbb{R}^{n_1 \times n_2}$, $\mathbf{A}_{21} \in \mathbb{R}^{n_2 \times n_1}$, $\mathbf{A}_{22} \in \mathbb{R}^{n_2 \times n_2}$, $\mathbf{B}_1 \in \mathbb{R}^{n_1 \times m}$, $\mathbf{B}_2 \in \mathbb{R}^{n_2 \times m}$, $\mathbf{C}_1 \in \mathbb{R}^{p \times n_1}$, $\mathbf{C}_2 \in \mathbb{R}^{p \times n_2}$, $\mathbf{D} \in \mathbb{R}^{p \times m}$. We assume that $\mathbf{A}_{22}$ and $\mathbf{E}_{11} - \mathbf{E}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21}$ are both nonsingular. In this case system (5.1) is of index 1. We now compute the polynomial part of this system.

PROPOSITION 5.1. *Let* $\mathbf{G}(s)$ *be a transfer function of the descriptor system* (5.1), *where* $\mathbf{A}_{22}$ *and* $\mathbf{E}_{11} - \mathbf{E}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21}$ *are both nonsingular. Then the polynomial part of* $\mathbf{G}(s)$ *is a constant matrix given by*

$$\mathbf{P}(s) = \mathbf{C}_1\mathbf{M}_1\mathbf{B}_2 + \mathbf{C}_2\mathbf{M}_2\mathbf{B}_2 + \mathbf{D},$$

*where*

$$\mathbf{M}_1 = (\mathbf{E}_{11} - \mathbf{E}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21})^{-1}\mathbf{E}_{12}\mathbf{A}_{22}^{-1}, \tag{5.2}$$
$$\mathbf{M}_2 = -\mathbf{A}_{22}^{-1}\mathbf{A}_{21}(\mathbf{E}_{11} - \mathbf{E}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21})^{-1}\mathbf{E}_{12}\mathbf{A}_{22}^{-1} - \mathbf{A}_{22}^{-1}. \tag{5.3}$$

*Proof.* Consider

$$(s\mathbf{E} - \mathbf{A})^{-1}\mathbf{B} = \begin{bmatrix} s\mathbf{E}_{11} - \mathbf{A}_{11} & s\mathbf{E}_{12} - \mathbf{A}_{12} \\ -\mathbf{A}_{21} & -\mathbf{A}_{22} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{B}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{F}_1(s) \\ \mathbf{F}_2(s) \end{bmatrix}.$$

This leads to

$$(s\mathbf{E}_{11} - \mathbf{A}_{11})\mathbf{F}_1(s) + (s\mathbf{E}_{12} - \mathbf{A}_{12})\mathbf{F}_2(s) = \mathbf{B}_1, \tag{5.4}$$
$$-\mathbf{A}_{21}\mathbf{F}_1(s) - \mathbf{A}_{22}\mathbf{F}_2(s) = \mathbf{B}_2. \tag{5.5}$$

Solving (5.5) for $\mathbf{F}_2(s)$ gives $\mathbf{F}_2(s) = -\mathbf{A}_{22}^{-1}(\mathbf{B}_2 + \mathbf{A}_{21}\mathbf{F}_1(s))$, and, thus,

$$\mathbf{F}_1(s) = \left((s\mathbf{E}_{11} - \mathbf{A}_{11}) - (s\mathbf{E}_{12} - \mathbf{A}_{12})\mathbf{A}_{22}^{-1}\mathbf{A}_{21}\right)^{-1}\left(\mathbf{B}_1 + (s\mathbf{E}_{12} - \mathbf{A}_{12})\mathbf{A}_{22}^{-1}\mathbf{B}_2\right)$$

implying that

$$\lim_{s\to\infty} \mathbf{F}_1(s) = \left(\mathbf{E}_{11} - \mathbf{E}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21}\right)^{-1}\mathbf{E}_{12}\mathbf{A}_{22}^{-1}\mathbf{B}_2.$$

Taking into account (5.5), we have

$$\lim_{s\to\infty} \mathbf{F}_2(s) = \left[-\mathbf{A}_{22}^{-1}\mathbf{A}_{21}\left(\mathbf{E}_{11} - \mathbf{E}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21}\right)^{-1}\mathbf{E}_{12}\mathbf{A}_{22}^{-1} - \mathbf{A}_{22}^{-1}\right]\mathbf{B}_2.$$

Finally, note that $\mathbf{P}(s) = \lim_{s\to\infty} \mathbf{G}(s) = \lim_{s\to\infty}(\mathbf{C}_1\mathbf{F}_1(s) + \mathbf{C}_2\mathbf{F}_2(s) + \mathbf{D})$, which leads to the desired conclusion. $\square$

We are now ready to state the interpolation result for the descriptor system (5.1). This result was briefly hinted at in the recent survey [1]. Here, we present it with a formal proof together with the formula developed for $\mathbf{P}(s)$ in Proposition 5.1. As our main focus will be $\mathcal{H}_2$-based model reduction, we will list the interpolation conditions only for the bi-tangential Hermite interpolation. Extension to the higher-order derivative interpolation is straightforward as shown in the earlier sections.

LEMMA 5.2. *Let* $\mathbf{G}(s)$ *be a transfer function of the semi-explicit descriptor system* (5.1). *For given* $r$ *distinct interpolation points* $\{\sigma_i\}_{i=1}^r$, *left tangential directions*

$\{c_i\}_{i=1}^r$ *and right tangential directions* $\{b_i\}_{i=1}^r$, *let* $\mathbf{V} \in \mathbb{C}^{n \times r}$ *and* $\mathbf{W} \in \mathbb{C}^{n \times r}$ *be given by*

$$\mathbf{V} = [\,(\sigma_1 \mathbf{E} - \mathbf{A})^{-1} \mathbf{B} b_1, \ \ldots, \ (\sigma_r \mathbf{E} - \mathbf{A})^{-1} \mathbf{B} b_r\,], \tag{5.6}$$

$$\mathbf{W} = [\,(\sigma_1 \mathbf{E} - \mathbf{A})^{-T} \mathbf{C}^T c_1, \ \ldots, \ (\sigma_r \mathbf{E} - \mathbf{A})^{-T} \mathbf{C}^T c_r\,]. \tag{5.7}$$

*Furthermore, let* $\mathcal{B}$ *and* $\mathcal{C}$ *be the matrices composed of the tangential directions as*

$$\mathcal{B} = [\,b_1, \ \ldots, \ b_r\,] \qquad and \qquad \mathcal{C} = [\,c_1, \ \ldots, \ c_r\,]. \tag{5.8}$$

*Define the reduced-order system matrices as*

$$\widetilde{\mathbf{E}} = \mathbf{W}^T \mathbf{E} \mathbf{V}, \qquad \widetilde{\mathbf{A}} = \mathbf{W}^T \mathbf{A} \mathbf{V} + \mathcal{C}^T \widetilde{\mathbf{D}} \mathcal{B}, \quad \widetilde{\mathbf{B}} = \mathbf{W}^T \mathbf{B} - \mathcal{C}^T \widetilde{\mathbf{D}},$$
$$\widetilde{\mathbf{C}} = \mathbf{C} \mathbf{V} - \widetilde{\mathbf{D}} \mathcal{B}, \qquad \widetilde{\mathbf{D}} = \mathbf{C}_1 \mathbf{M}_1 \mathbf{B}_2 + \mathbf{C}_2 \mathbf{M}_2 \mathbf{B}_2 + \mathbf{D}. \tag{5.9}$$

*Then the polynomial parts of* $\widetilde{\mathbf{G}}(s) = \widetilde{\mathbf{C}}(s\widetilde{\mathbf{E}} - \widetilde{\mathbf{A}})^{-1}\widetilde{\mathbf{B}} + \widetilde{\mathbf{D}}$ *and* $\mathbf{G}(s)$ *match assuming* $\widetilde{\mathbf{E}}$ *is nonsingular, and* $\widetilde{\mathbf{G}}(s)$ *satisfies the bi-tangential Hermite interpolation conditions*

$$\mathbf{G}(\sigma_i) b_i = \widetilde{\mathbf{G}}(\sigma_i) b_i, \qquad c_i^T \mathbf{G}(\sigma_i) = c_i^T \widetilde{\mathbf{G}}(\sigma_i), \qquad c_i^T \mathbf{G}'(\sigma_i) b_i = c_i^T \widetilde{\mathbf{G}}'(\sigma_i) b_i$$

*for* $i = 1, \ldots, r$, *provided* $\sigma_i \mathbf{E} - \mathbf{A}$ *and* $\sigma_i \widetilde{\mathbf{E}} - \widetilde{\mathbf{A}}$ *are both nonsingular.*

*Proof.* Since $\widetilde{\mathbf{E}}$ is nonsingular, $\lim\limits_{s \to \infty} \widetilde{\mathbf{G}}(s) = \widetilde{\mathbf{D}}$. But by Lemma 5.1, we have $\widetilde{\mathbf{D}} = \lim\limits_{s \to \infty} \mathbf{G}(s)$ ensuring that the polynomial parts of $\mathbf{G}(s)$ and $\widetilde{\mathbf{G}}(s)$ coincide. The interpolation property is a result of [2, 22], where it is shown that the appropriate shifting of the reduced-order quantities with a non-zero feedthrough term as done in (5.9) attains the original bi-tangential interpolation conditions hidden in $\mathbf{V}$ and $\mathbf{W}$ of (5.6) and (5.7), respectively. $\square$

This result leads to Algorithm 5.1, which achieves bi-tangential Hermite interpolation of the semi-explicit descriptor system (5.1) without explicitly forming the spectral projectors.

---

ALGORITHM 5.1. **Interpolatory model reduction for semi-explicit descriptor systems of index 1**

1) *Make an initial selection of the interpolation points* $\{\sigma_i\}_{i=1}^r$ *and the tangential directions* $\{b_i\}_{i=1}^r$ *and* $\{c_i\}_{i=1}^r$.
2) $\mathbf{V} = [\,(\sigma_1 \mathbf{E} - \mathbf{A})^{-1} \mathbf{B} b_1, \ \ldots, \ (\sigma_r \mathbf{E} - \mathbf{A})^{-1} \mathbf{B} b_r\,]$,

 $\mathbf{W} = [\,(\sigma_1 \mathbf{E} - \mathbf{A})^{-T} \mathbf{C}^T c_1, \ \ldots, \ (\sigma_r \mathbf{E} - \mathbf{A})^{-T} \mathbf{C}^T c_r\,]$.
3) *Define* $\widetilde{\mathbf{D}} = \mathbf{C}_1 \mathbf{M}_1 \mathbf{B}_2 + \mathbf{C}_2 \mathbf{M}_2 \mathbf{B}_2 + \mathbf{D}$, *where* $\mathbf{M}_1$ *and* $\mathbf{M}_2$ *are defined in (5.2) and (5.3), respectively.*
4) *Define* $\mathcal{B} = [\,b_1, \ \ldots, \ b_r\,]$ *and* $\mathcal{C} = [\,c_1, \ \ldots, \ c_r\,]$.
5) $\widetilde{\mathbf{E}} = \mathbf{W}^T \mathbf{E} \mathbf{V}, \ \ \widetilde{\mathbf{A}} = \mathbf{W}^T \mathbf{A} \mathbf{V} + \mathcal{C}^T \widetilde{\mathbf{D}} \mathcal{B}, \ \ \widetilde{\mathbf{B}} = \mathbf{W}^T \mathbf{B} - \mathcal{C}^T \widetilde{\mathbf{D}}, \ \ \widetilde{\mathbf{C}} = \mathbf{C} \mathbf{V} - \widetilde{\mathbf{D}} \mathcal{B}$.

---

We want to emphasize that the assumption in Lemma 5.2 that $\widetilde{\mathbf{E}}$ be nonsingular is not restrictive. This will be the case generically. If $\mathbf{W}$ and $\mathbf{V}$ are full-rank $n \times r$ matrices, the rank of the $r \times r$ matrix $\widetilde{\mathbf{E}} = \mathbf{W}^T \mathbf{E} \mathbf{V}$ will be generically $r$ as long as $rank(\mathbf{E}) > r$. The fact that $\widetilde{\mathbf{E}}$ will be full-rank is indeed the precise reason why we cannot simply apply Theorem 2.1 to descriptor systems.

**5.1. Optimal $\mathcal{H}_2$ model reduction for semi-explicit descriptor systems.** Lemma 5.2 provides the theoretical basis for an IRKA-based iteration for $\mathcal{H}_2$ model reduction of semi-explicit descriptor systems. One naive approach would be the following: Given system (5.1), simply apply IRKA of [15] to obtain an intermediate reduced-order model $\widehat{\mathbf{G}}(s) = \widehat{\mathbf{C}}(s\widehat{\mathbf{E}} - \widehat{\mathbf{A}})^{-1}\widehat{\mathbf{B}} + \widehat{\mathbf{D}}$. Of course, this will be generically an ODE and will not necessarily match the behavior of $\mathbf{G}(s)$ around $s = \infty$. Thus, apply Lemma 5.2 to obtain the final reduced-model $\widetilde{\mathbf{G}}(s) = \widetilde{\mathbf{C}}(s\widetilde{\mathbf{E}} - \widetilde{\mathbf{A}})^{-1}\widetilde{\mathbf{B}} + \widetilde{\mathbf{D}}$ with

$$\widetilde{\mathbf{E}} = \widehat{\mathbf{E}}, \quad \widetilde{\mathbf{A}} = \widehat{\mathbf{A}} + \mathcal{C}^T\widetilde{\mathbf{D}}\mathcal{B}, \quad \widetilde{\mathbf{B}} = \widehat{\mathbf{B}} - \mathcal{C}^T\widetilde{\mathbf{D}}, \quad \widetilde{\mathbf{C}} = \widehat{\mathbf{C}} - \widetilde{\mathbf{D}}\mathcal{B}, \qquad (5.10)$$

where $\widetilde{\mathbf{D}}$ is defined as in Lemma 5.2. While this shifting of the intermediate matrices by the $\mathbf{D}$-term guarantees that the polynomial parts of $\mathbf{G}(s)$ and $\widetilde{\mathbf{G}}(s)$ match, the $\mathcal{H}_2$ optimality conditions will *not* be satisfied. The reason is as follows. Recall that the $\mathcal{H}_2$ optimality requires bi-tangential Hermite interpolation at the mirror images of the reduced-order poles. The intermediate model $\widehat{\mathbf{G}}(s)$ satisfies this but since it does not match the polynomial part, the resulting $\mathcal{H}_2$ error is unbounded. Then constructing $\widetilde{\mathbf{G}}(s)$ as in (5.10), we enforce the matching of the polynomial part but $\widetilde{\mathbf{G}}(s)$ still interpolates $\mathbf{G}(s)$ at the same interpolation points as $\widehat{\mathbf{G}}(s)$, i.e., at the mirror images of the poles of $\widehat{\mathbf{G}}(s)$. However, clearly due to (5.10), the poles of $\widehat{\mathbf{G}}(s)$ and $\widetilde{\mathbf{G}}(s)$ are different; thus $\widetilde{\mathbf{G}}(s)$ will no longer satisfy the optimal $\mathcal{H}_2$ necessary conditions. In order to achieve both the mirror-image interpolation conditions and the polynomial part matching, the $\widetilde{\mathbf{D}}$ term modification must be included throughout the iteration, not just at the end. This results in Algorithm 5.2.

---

ALGORITHM 5.2. **IRKA for semi-explicit descriptor systems of index 1**
1) *Make an initial shift selection $\{\sigma_i\}_{i=1}^r$ and initial tangential directions $\{\mathsf{b}_i\}_{i=1}^r$ and $\{\mathsf{c}_i\}_{i=1}^r$.*
2) $\mathbf{V}_r = \left[(\sigma_1\mathbf{E} - \mathbf{A})^{-1}\mathbf{B}\mathsf{b}_1, \ldots, (\sigma_r\mathbf{E} - \mathbf{A})^{-1}\mathbf{B}\mathsf{b}_r\right]$,
   $\mathbf{W}_r = \left[(\sigma_1\mathbf{E} - \mathbf{A})^{-T}\mathbf{C}^T\mathsf{c}_1, \ldots, (\sigma_r\mathbf{E} - \mathbf{A})^{-T}\mathbf{C}^T\mathsf{c}_r\right]$.
3) *Define $\widetilde{\mathbf{D}} = \mathbf{C}_1\mathbf{M}_1\mathbf{B}_2 + \mathbf{C}_2\mathbf{M}_2\mathbf{B}_2 + \mathbf{D}$, where $\mathbf{M}_1$ and $\mathbf{M}_2$ are defined in (5.2) and (5.3), respectively.*
4) *Define $\mathcal{B} = [\mathsf{b}_1, \ldots, \mathsf{b}_r]$ and $\mathcal{C} = [\mathsf{c}_1, \ldots, \mathsf{c}_r]$.*
5) *while (not converged)*
   a) $\widetilde{\mathbf{E}} = \mathbf{W}^T\mathbf{E}\mathbf{V}, \widetilde{\mathbf{A}} = \mathbf{W}^T\mathbf{A}\mathbf{V} + \mathcal{C}^T\widetilde{\mathbf{D}}\mathcal{B}, \widetilde{\mathbf{B}} = \mathbf{W}^T\mathbf{B} - \mathcal{C}^T\widetilde{\mathbf{D}}, \widetilde{\mathbf{C}} = \mathbf{C}\mathbf{V} - \widetilde{\mathbf{D}}\mathcal{B}$.
   b) *Compute $\mathbf{Y}^*\widetilde{\mathbf{A}}\mathbf{Z} = \text{diag}(\lambda_1, \ldots, \lambda_r)$ and $\mathbf{Y}^*\widetilde{\mathbf{E}}\mathbf{Z} = \mathbf{I}_r$, where the columns of $\mathbf{Z} = [\mathbf{z}_1, \ldots, \mathbf{z}_r]$ and $\mathbf{Y} = [\mathbf{y}_1, \ldots, \mathbf{y}_r]$ are, respectively, the right and left eigenvectors of $\lambda\widetilde{\mathbf{E}} - \widetilde{\mathbf{A}}$.*
   c) $\sigma_i \leftarrow -\lambda_i, \mathsf{b}_i^T \leftarrow \mathbf{y}_i^*\widetilde{\mathbf{B}}$ and $\mathsf{c}_i \leftarrow \widetilde{\mathbf{C}}\mathbf{z}_i$ *for $i = 1, \ldots, r$.*
   d) $\mathbf{V} = \left[(\sigma_1\mathbf{E} - \mathbf{A})^{-1}\mathbf{B}\mathsf{b}_1, \ldots, (\sigma_r\mathbf{E} - \mathbf{A})^{-1}\mathbf{B}\mathsf{b}_r\right]$,
      $\mathbf{W} = \left[(\sigma_1\mathbf{E} - \mathbf{A})^{-T}\mathbf{C}^T\mathsf{c}_1, \ldots, (\sigma_r\mathbf{E} - \mathbf{A})^{-T}\mathbf{C}^T\mathsf{c}_r\right]$.
   *end while*
6) $\widetilde{\mathbf{E}} = \mathbf{W}^T\mathbf{E}\mathbf{V}, \widetilde{\mathbf{A}} = \mathbf{W}^T\mathbf{A}\mathbf{V} + \mathcal{C}^T\widetilde{\mathbf{D}}\mathcal{B}, \widetilde{\mathbf{B}} = \mathbf{W}^T\mathbf{B} - \mathcal{C}^T\widetilde{\mathbf{D}}, \widetilde{\mathbf{C}} = \mathbf{C}\mathbf{V} - \widetilde{\mathbf{D}}\mathcal{B}$.

---

The next result is a restatement of the above discussion.

COROLLARY 5.3. *Let $\mathbf{G}(s)$ be a transfer function of the semi-explicit descriptor system (5.1) and let $\widetilde{\mathbf{G}}(s) = \widetilde{\mathbf{C}}(s\widetilde{\mathbf{E}} - \widetilde{\mathbf{A}})^{-1}\widetilde{\mathbf{B}} + \widetilde{\mathbf{D}}$ be obtained by Algorithm 5.2. Then $\widetilde{\mathbf{G}}(s)$ satisfies the first-order necessary conditions of the $\mathcal{H}_2$ optimal model reduction problem.*

**5.2. Supersonic inlet flow example.** Consider the Euler equations modelling the unsteady flow through a supersonic diffuser as described in [21]. Linearization around a steady-state solution and spatial discretization using a finite volume method leads to a semi-explicit descriptor system (5.1) of dimension $n = 11730$. For simplicity, we focus on the single-input single-output subsystem dynamics corresponding to the input as the bleed actuation mass flow and the output as the average Mach number.

It is important to emphasize that applying balanced truncation to this model is far from trivial because of difficulty of solving the Lyapunov equations. Instead, we apply the proposed method in Algorithm 5.2 to obtain an $\mathcal{H}_2$-optimal reduced-model of order $r = 11$, where the only cost are sparse linear solves and the need for computing the spectral projectors are removed. As pointed out in [21], the frequencies of practical interest are the low frequency components. Figure 5.2 shows the amplitude and phase plots of $\mathbf{G}(\imath\omega)$ and $\widetilde{\mathbf{G}}(\imath\omega)$ for $\omega \in [0, 25]$ illustrating a very accurate match of the original model. The resulting model reduction errors are

$$\frac{\|\mathbf{G} - \widetilde{\mathbf{G}}\|_{\mathcal{H}_\infty}}{\|\mathbf{G}\|_{\mathcal{H}_\infty}} = 5.2252 \times 10^{-2} \qquad \text{and} \qquad \frac{\|\mathbf{G}_{\mathrm{sp}} - \widetilde{\mathbf{G}}_{\mathrm{sp}}\|_{\mathcal{H}_\infty}}{\|\mathbf{G}_{\mathrm{sp}}\|_{\mathcal{H}_\infty}} = 5.2251 \times 10^{-2}.$$
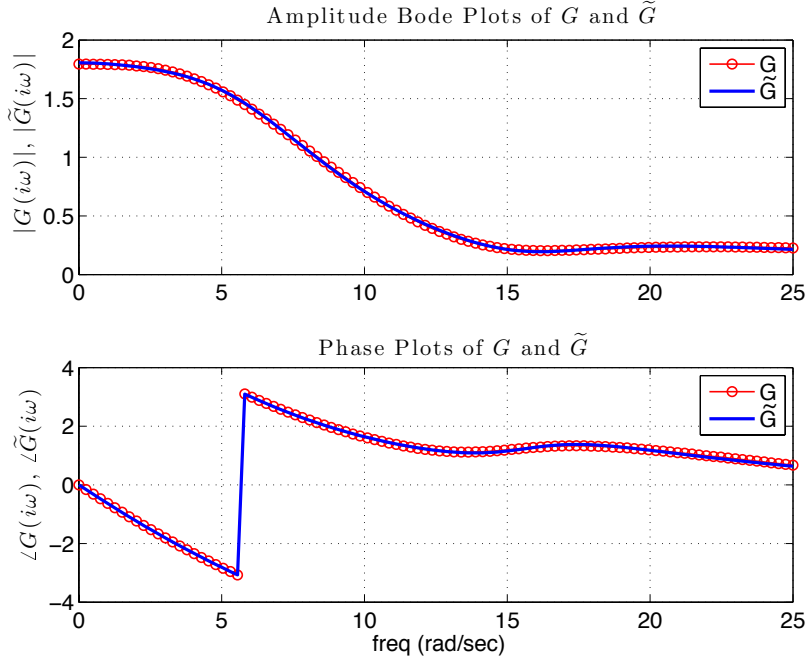


FIG. 5.1. *Supersonic inlet flow model: amplitude and phase Bode plots of* $\mathbf{G}(s)$ *and* $\widetilde{\mathbf{G}}(s)$.

**6. Stokes-type descriptor systems of index 2.** In this section, we consider a Stokes-type descriptor system of the form

$$\begin{aligned} \mathbf{E}_{11}\dot{\mathbf{x}}_1(t) &= \mathbf{A}_{11}\mathbf{x}_1(t) + \mathbf{A}_{12}\mathbf{x}_2(t) + \mathbf{B}_1\mathbf{u}(t), \\ \mathbf{0} &= \mathbf{A}_{21}\mathbf{x}_1(t) + \mathbf{B}_2\mathbf{u}(t), \\ \mathbf{y}(t) &= \mathbf{C}_1\mathbf{x}_1(t) + \mathbf{C}_2\mathbf{x}_2(t) + \mathbf{D}\mathbf{u}(t), \end{aligned} \qquad (6.1)$$

where the state is $\mathbf{x}(t) = \left[\mathbf{x}_1^T(t), \mathbf{x}_2^T(t)\right]^T \in \mathbb{R}^n$ with $\mathbf{x}_1(t) \in \mathbb{R}^{n_1}$, $\mathbf{x}_2(t) \in \mathbb{R}^{n_2}$ and $n_1 + n_2 = n$, the input is $\mathbf{u}(t) \in \mathbb{R}^m$, the output is $\mathbf{y}(t) \in \mathbb{R}^p$, and $\mathbf{E}_{11}, \mathbf{A}_{11} \in \mathbb{R}^{n_1 \times n_1}$,

$\mathbf{A}_{12} \in \mathbb{R}^{n_1 \times n_2}$, $\mathbf{A}_{21} \in \mathbb{R}^{n_2 \times n_1}$, $\mathbf{B}_1 \in \mathbb{R}^{n_1 \times m}$, $\mathbf{B}_2 \in \mathbb{R}^{n_2 \times m}$, $\mathbf{C}_1 \in \mathbb{R}^{p \times n_1}$, $\mathbf{C}_2 \in \mathbb{R}^{p \times n_2}$, and $\mathbf{D} \in \mathbb{R}^{p \times m}$. We assume that $\mathbf{E}_{11}$ is nonsingular, $\mathbf{A}_{12}$ and $\mathbf{A}_{21}^T$ have both full column rank and $\mathbf{A}_{21}\mathbf{E}_{11}^{-1}\mathbf{A}_{12}$ is nonsingular. In this case, system (6.1) is of index 2.

In [17], the authors showed how to apply ADI-based balanced truncation to systems of the form (6.1) without explicit projector computation. Here, we extend this analysis to interpolatory model reduction and show how to reduce (6.1) optimally in the $\mathcal{H}_2$-norm without computing the deflating subspaces. Unlike [17], $\mathbf{E}_{11}$ is not assumed to be symmetric and positive definite, and $\mathbf{A}_{21}$ is not assumed to be equal to $\mathbf{A}_{12}^T$.

First, consider system (6.1) with $\mathbf{B}_2 = \mathbf{0}$, as the case of $\mathbf{B}_2 \neq \mathbf{0}$ follows similarly. Following the exposition of [17], consider the projectors

$$\mathbf{\Pi}_l = \mathbf{I} - \mathbf{E}_{11}^{-1}\mathbf{A}_{12}(\mathbf{A}_{21}\mathbf{E}_{11}^{-1}\mathbf{A}_{12})^{-1}\mathbf{A}_{21},$$
$$\mathbf{\Pi}_r = \mathbf{I} - \mathbf{A}_{12}(\mathbf{A}_{21}\mathbf{E}_{11}^{-1}\mathbf{A}_{12})^{-1}\mathbf{A}_{21}\mathbf{E}_{11}^{-1}.$$

Then the descriptor system (6.1) can be decoupled into a system

$$\begin{aligned} \mathbf{\Pi}_l\mathbf{E}_{11}\mathbf{\Pi}_r\dot{\mathbf{x}}_1(t) &= \mathbf{\Pi}_l\mathbf{A}_{11}\mathbf{\Pi}_r\mathbf{x}_1(t) + \mathbf{\Pi}_l\mathbf{B}_1\mathbf{u}(t) \\ \mathbf{y}(t) &= \mathbf{C}\mathbf{\Pi}_r\mathbf{x}_1(t) + \boldsymbol{\mathcal{D}}\mathbf{u}(t) \end{aligned} \tag{6.2}$$

with

$$\begin{aligned} \mathbf{C} &= \mathbf{C}_1 - \mathbf{C}_2(\mathbf{A}_{21}\mathbf{E}_{11}^{-1}\mathbf{A}_{12})^{-1}\mathbf{A}_{21}\mathbf{E}_{11}^{-1}\mathbf{A}_{11}, \\ \boldsymbol{\mathcal{D}} &= \mathbf{D} - \mathbf{C}_2(\mathbf{A}_{21}\mathbf{E}_{11}^{-1}\mathbf{A}_{12})^{-1}\mathbf{A}_{21}\mathbf{E}_{11}^{-1}\mathbf{B}_1, \end{aligned}$$

and an algebraic equation

$$\mathbf{x}_2(t) = -(\mathbf{A}_{21}\mathbf{E}_{11}^{-1}\mathbf{A}_{12})^{-1}\mathbf{A}_{21}\mathbf{E}_{11}^{-1}\mathbf{A}_{11}\mathbf{x}_1(t) - (\mathbf{A}_{21}\mathbf{E}_{11}^{-1}\mathbf{A}_{12})^{-1}\mathbf{A}_{21}\mathbf{E}_{11}^{-1}\mathbf{B}_1\mathbf{u}(t).$$

By decomposing $\mathbf{\Pi}_l$ and $\mathbf{\Pi}_r$ as

$$\mathbf{\Pi}_l = \mathbf{\Theta}_{l,1}\mathbf{\Theta}_{l,2}^T, \qquad \mathbf{\Pi}_r = \mathbf{\Theta}_{r,1}\mathbf{\Theta}_{r,2}^T \tag{6.3}$$

with $\mathbf{\Theta}_{l,j}, \mathbf{\Theta}_{r,j} \in \mathbb{R}^{n_1 \times (n_1 - n_2)}$ such that

$$\mathbf{\Theta}_{l,2}^T\mathbf{\Theta}_{l,1} = \mathbf{I}, \qquad \mathbf{\Theta}_{r,2}^T\mathbf{\Theta}_{r,1} = \mathbf{I}, \tag{6.4}$$

and defining $\tilde{\mathbf{x}}_1(t) = \mathbf{\Theta}_{r,2}^T\mathbf{x}_1(t)$, system (6.2) becomes

$$\begin{aligned} \mathbf{\Theta}_{l,2}^T\mathbf{E}_{11}\mathbf{\Theta}_{r,1}\dot{\tilde{\mathbf{x}}}_1(t) &= \mathbf{\Theta}_{l,2}^T\mathbf{A}_{11}\mathbf{\Theta}_{r,1}\tilde{\mathbf{x}}_1(t) + \mathbf{\Theta}_{l,2}^T\mathbf{B}_1\mathbf{u}(t) \\ \mathbf{y}(t) &= \mathbf{C}\mathbf{\Theta}_{r,1}\tilde{\mathbf{x}}_1(t) + \boldsymbol{\mathcal{D}}\mathbf{u}(t). \end{aligned} \tag{6.5}$$

Then the reduction of the descriptor system (6.1) is equivalent to the reduction of system (6.2) or (6.5). However, the beauty of this equivalence lies in the observation that the matrix $\mathbf{\Theta}_{l,2}^T\mathbf{E}_{11}\mathbf{\Theta}_{r,1}$ is nonsingular. Therefore, standard model reduction procedures for ODEs can be applied to system (6.5), and the obtained reduced-order model will approximate the descriptor system (6.1). It is important to emphasize that even though (6.2) and (6.5) are equivalent to (6.1), the ultimate goal of this section is to develop an interpolatory model reduction method that does not require the explicit computation of either the projectors $\mathbf{\Pi}_l$, $\mathbf{\Pi}_r$ or the basis matrices $\mathbf{\Theta}_{l,2}$, $\mathbf{\Theta}_{r,1}$. For this purpose, define the matrices

$$\boldsymbol{\mathcal{E}} = \mathbf{\Pi}_l\mathbf{E}_{11}\mathbf{\Pi}_r, \qquad \boldsymbol{\mathcal{A}} = \mathbf{\Pi}_l\mathbf{A}_{11}\mathbf{\Pi}_r, \qquad \boldsymbol{\mathcal{B}} = \mathbf{\Pi}_l\mathbf{B}_1, \quad \boldsymbol{\mathcal{C}} = \mathbf{C}\mathbf{\Pi}_r. \tag{6.6}$$

In interpolation setting, the matrix of interest will be $\sigma\boldsymbol{\mathcal{E}} - \boldsymbol{\mathcal{A}}$ with $\sigma \in \mathbb{C}$. Luckily, several key properties of $\boldsymbol{\mathcal{E}} + \tau\boldsymbol{\mathcal{A}}$, $\tau \in \mathbb{C}$, were already introduced in [17]. However, we present these results in terms of $\sigma\boldsymbol{\mathcal{E}} - \boldsymbol{\mathcal{A}}$ instead of $\boldsymbol{\mathcal{E}} + \tau\boldsymbol{\mathcal{A}}$.

LEMMA 6.1. *Let $\boldsymbol{\Theta}_{l,2}$ and $\boldsymbol{\Theta}_{r,1}$ be the matrices defined in (6.3) and let $\sigma \in \mathbb{C}$ be such that $\sigma\boldsymbol{\Theta}_{l,2}^T\mathbf{E}_{11}\boldsymbol{\Theta}_{r,1} - \boldsymbol{\Theta}_{l,2}^T\mathbf{A}_{11}\boldsymbol{\Theta}_{r,1}$ is nonsingular. The matrix defined as*

$$(\sigma\boldsymbol{\mathcal{E}} - \boldsymbol{\mathcal{A}})^I := \boldsymbol{\Theta}_{r,1}(\sigma\boldsymbol{\Theta}_{l,2}^T\mathbf{E}_{11}\boldsymbol{\Theta}_{r,1} - \boldsymbol{\Theta}_{l,2}^T\mathbf{A}_{11}\boldsymbol{\Theta}_{r,1})^{-1}\boldsymbol{\Theta}_{l,2}^T \tag{6.7}$$

*satisfies*

$$(\sigma\boldsymbol{\mathcal{E}} - \boldsymbol{\mathcal{A}})^I(\sigma\boldsymbol{\mathcal{E}} - \boldsymbol{\mathcal{A}}) = \boldsymbol{\Pi}_r \quad and \quad (\sigma\boldsymbol{\mathcal{E}} - \boldsymbol{\mathcal{A}})(\sigma\boldsymbol{\mathcal{E}} - \boldsymbol{\mathcal{A}})^I = \boldsymbol{\Pi}_l.$$

*Similarly, the matrix defined as*

$$(\sigma\boldsymbol{\mathcal{E}}^T - \boldsymbol{\mathcal{A}}^T)^I := \boldsymbol{\Theta}_{l,2}(\sigma\boldsymbol{\Theta}_{r,1}^T\mathbf{E}_{11}^T\boldsymbol{\Theta}_{l,2} - \boldsymbol{\Theta}_{r,1}^T\mathbf{A}_{11}^T\boldsymbol{\Theta}_{l,2})^{-1}\boldsymbol{\Theta}_{r,1}^T \tag{6.8}$$

*satisfies*

$$(\sigma\boldsymbol{\mathcal{E}}^T - \boldsymbol{\mathcal{A}}^T)^I(\sigma\boldsymbol{\mathcal{E}}^T - \boldsymbol{\mathcal{A}}^T) = \boldsymbol{\Pi}_l^T \quad and \quad (\sigma\boldsymbol{\mathcal{E}}^T - \boldsymbol{\mathcal{A}}^T)(\sigma\boldsymbol{\mathcal{E}}^T - \boldsymbol{\mathcal{A}}^T)^I = \boldsymbol{\Pi}_r^T.$$

*Proof.* Following a similar argument to that in [17], the proof of the first equality follows directly from (6.3) and (6.7). Indeed, we have

$$\begin{aligned}(\sigma\boldsymbol{\mathcal{E}}-\boldsymbol{\mathcal{A}})^I(\sigma\boldsymbol{\mathcal{E}} - \boldsymbol{\mathcal{A}}) &= \boldsymbol{\Theta}_{r,1}(\sigma\boldsymbol{\Theta}_{l,2}^T\mathbf{E}_{11}\boldsymbol{\Theta}_{r,1} - \boldsymbol{\Theta}_{l,2}^T\mathbf{A}_{11}\boldsymbol{\Theta}_{r,1})^{-1}\boldsymbol{\Theta}_{l,2}^T\boldsymbol{\Pi}_l(\sigma\mathbf{E}_{11} - \mathbf{A}_{11})\boldsymbol{\Pi}_r \\ &= \boldsymbol{\Theta}_{r,1}(\sigma\boldsymbol{\Theta}_{l,2}^T\mathbf{E}_{11}\boldsymbol{\Theta}_{r,1} - \boldsymbol{\Theta}_{l,2}^T\mathbf{A}_{11}\boldsymbol{\Theta}_{r,1})^{-1}\boldsymbol{\Theta}_{l,2}^T(\sigma\mathbf{E}_{11} - \mathbf{A}_{11})\boldsymbol{\Theta}_{r,1}\boldsymbol{\Theta}_{r,2}^T \\ &= \boldsymbol{\Theta}_{r,1}\boldsymbol{\Theta}_{r,2}^T = \boldsymbol{\Pi}_r.\end{aligned}$$

The remaining equalities follow similarly. $\square$

At first glance, the definition of the generalized inverses in (6.7) and (6.8) may seem to be irrelevant for model reduction of the descriptor system (6.1). Recall that reducing (6.1) is equivalent to reducing system (6.2) and the interpolatory projection method for (6.2) will require inverting $(\sigma\boldsymbol{\mathcal{E}} - \boldsymbol{\mathcal{A}})$ and $(\sigma\boldsymbol{\mathcal{E}}^T - \boldsymbol{\mathcal{A}}^T)$. However, these inverses do not exist. As a result, definitions (6.7) and (6.8) become pivotal in order to achieve interpolatory model reduction of (6.2) and, thereby, of (6.1) as shown in the next theorem.

THEOREM 6.2. *Let $s = \sigma, \mu \in \mathbb{C}$ be such that the matrices*

$$s\boldsymbol{\Theta}_{l,2}^T\mathbf{E}_{11}\boldsymbol{\Theta}_{r,1} - \boldsymbol{\Theta}_{l,2}^T\mathbf{A}_{11}\boldsymbol{\Theta}_{r,1} \qquad and \qquad s\mathbf{W}^T\mathbf{E}_{11}\mathbf{V} - \mathbf{W}^T\mathbf{A}_{11}\mathbf{V}$$

*are invertible. Define the reduced-order model*

$$\widetilde{\mathbf{G}}(s) = \mathbf{CV}(s\mathbf{W}^T\mathbf{E}_{11}\mathbf{V} - \mathbf{W}^T\mathbf{A}_{11}\mathbf{V})^{-1}\mathbf{W}^T\mathbf{B}_1 + \boldsymbol{\mathcal{D}}. \tag{6.9}$$

*Let $\mathsf{b} \in \mathbb{C}^m$ and $\mathsf{c} \in \mathbb{C}^p$ be fixed nontrivial vectors.*
1. *If $(\sigma\boldsymbol{\mathcal{E}} - \boldsymbol{\mathcal{A}})^I\boldsymbol{\mathcal{B}}\mathsf{b} \in \mathrm{Im}(\mathbf{V}) \subset \mathrm{Im}(\boldsymbol{\Theta}_{r,1})$ and $(\mu\boldsymbol{\mathcal{E}}^T - \boldsymbol{\mathcal{A}}^T)^I\boldsymbol{\mathcal{C}}^T\mathsf{c} \in \mathrm{Im}(\mathbf{W}) \subset \mathrm{Im}(\boldsymbol{\Theta}_{l,2})$, then $\mathbf{G}(\sigma)\mathsf{b} = \widetilde{\mathbf{G}}(\sigma)\mathsf{b}$ and $\mathsf{c}^T\mathbf{G}(\mu) = \mathsf{c}^T\widetilde{\mathbf{G}}(\mu)$.*
2. *If, in addition, $\sigma = \mu$, then $\mathsf{c}^T\mathbf{G}'(\sigma)\mathsf{b} = \mathsf{c}^T\widetilde{\mathbf{G}}'(\sigma)\mathsf{b}$.*

REMARK 6.1. *Before presenting the proof, we want to emphasize that this interpolation result is different than the usual interpolation framework given in Theorem 2.1,*

*where the projection matrices* $\mathbf{V}$ *and* $\mathbf{W}$ *are constructed using* $\mathbf{A}$, $\mathbf{E}$, $\mathbf{B}$ *and* $\mathbf{C}$ *and then the projection is applied to the same quantities. In Theorem 6.2, however, the projection matrices* $\mathbf{V}$ *and* $\mathbf{W}$ *are constructed using the system matrices of* (6.2), *namely* $\boldsymbol{\mathcal{A}}, \boldsymbol{\mathcal{E}}, \boldsymbol{\mathcal{B}}$ *and* $\boldsymbol{\mathcal{C}}$. *But then the projection (model reduction) is applied to the system matrices of* (6.1), *namely* $\mathbf{E}_{11}, \mathbf{A}_{11}, \mathbf{B}_1$ *and* $\mathbf{C}$. *Thus, the proof will serve to fill in this important gap.*

*Proof.* Since systems (6.1) and (6.5) are equivalent, they have the same transfer function given by

$$\mathbf{G}(s) = \mathbf{C}\boldsymbol{\Theta}_{r,1}(s\boldsymbol{\Theta}_{l,2}^T\mathbf{E}_{11}\boldsymbol{\Theta}_{r,1} - \boldsymbol{\Theta}_{l,2}^T\mathbf{A}_{11}\boldsymbol{\Theta}_{r,1})^{-1}\boldsymbol{\Theta}_{l,2}^T\mathbf{B}_1 + \boldsymbol{\mathcal{D}}.$$

Since $\boldsymbol{\Theta}_{l,2}^T\mathbf{E}_{11}\boldsymbol{\Theta}_{r,1}$ in (6.5) is nonsingular, we make use of Theorem 2.1. Define $\widetilde{\mathbf{V}}$ and $\widetilde{\mathbf{W}}$ such that

$$\mathbf{V} = \boldsymbol{\Theta}_{r,1}\widetilde{\mathbf{V}} \qquad \text{and} \qquad \mathbf{W} = \boldsymbol{\Theta}_{l,2}\widetilde{\mathbf{W}}. \tag{6.10}$$

Pluging these matrices into (6.9), we obtain that

$$\widetilde{\mathbf{G}}(s) = \mathbf{C}\boldsymbol{\Theta}_{r,1}\widetilde{\mathbf{V}}(s\widetilde{\mathbf{W}}^T\boldsymbol{\Theta}_{l,2}^T\mathbf{E}_{11}\boldsymbol{\Theta}_{r,1}\widetilde{\mathbf{V}} - \widetilde{\mathbf{W}}^T\boldsymbol{\Theta}_{l,2}^T\mathbf{A}_{11}\boldsymbol{\Theta}_{r,1}\widetilde{\mathbf{V}})^{-1}\widetilde{\mathbf{W}}^T\boldsymbol{\Theta}_{l,2}^T\mathbf{B}_1 + \boldsymbol{\mathcal{D}}.$$

Moreover, it follows from (6.4) that $\widetilde{\mathbf{V}} = \boldsymbol{\Theta}_{r,2}^T\mathbf{V}$ and $\widetilde{\mathbf{W}} = \boldsymbol{\Theta}_{l,1}^T\mathbf{W}$. To prove the first claim in part 1, we note that (6.3) implies that

$$\boldsymbol{\Theta}_{l,2}^T\boldsymbol{\mathcal{B}} = \boldsymbol{\Theta}_{l,2}^T\boldsymbol{\Pi}_l\mathbf{B}_1 = \boldsymbol{\Theta}_{l,2}^T\boldsymbol{\Theta}_{l,1}\boldsymbol{\Theta}_{l,2}^T\mathbf{B}_1 = \boldsymbol{\Theta}_{l,2}^T\mathbf{B}_1. \tag{6.11}$$

Since $(\sigma\boldsymbol{\mathcal{E}} - \boldsymbol{\mathcal{A}})^I\boldsymbol{\mathcal{B}}\mathbf{b} \in \text{Im}(\mathbf{V})$, there exists $\mathbf{q} \in \mathbb{R}^r$ such that $(\sigma\boldsymbol{\mathcal{E}} - \boldsymbol{\mathcal{A}})^I\boldsymbol{\mathcal{B}}\mathbf{b} = \mathbf{V}\mathbf{q}$. Using (6.7) (6.10) and (6.11), this equation can be written as

$$\boldsymbol{\Theta}_{r,1}(\sigma\boldsymbol{\Theta}_{l,2}^T\mathbf{E}_{11}\boldsymbol{\Theta}_{r,1} - \boldsymbol{\Theta}_{l,2}^T\mathbf{A}_{11}\boldsymbol{\Theta}_{r,1})^{-1}\boldsymbol{\Theta}_{l,2}^T\mathbf{B}_1\mathbf{b} = \boldsymbol{\Theta}_{r,1}\widetilde{\mathbf{V}}\mathbf{q}.$$

The left multiplication by $\boldsymbol{\Theta}_{r,2}^T$ gives

$$(\sigma\boldsymbol{\Theta}_{l,2}^T\mathbf{E}_{11}\boldsymbol{\Theta}_{r,1} - \boldsymbol{\Theta}_{l,2}^T\mathbf{A}_{11}\boldsymbol{\Theta}_{r,1})^{-1}\boldsymbol{\Theta}_{l,2}^T\mathbf{B}_1\mathbf{b} = \widetilde{\mathbf{V}}\mathbf{q}.$$

Hence, $(\sigma\boldsymbol{\Theta}_{l,2}^T\mathbf{E}_{11}\boldsymbol{\Theta}_{r,1} - \boldsymbol{\Theta}_{l,2}^T\mathbf{A}_{11}\boldsymbol{\Theta}_{r,1})^{-1}\boldsymbol{\Theta}_{l,2}^T\mathbf{B}_1\mathbf{b} \in \text{Im}(\widetilde{\mathbf{V}})$. Then it follows from Theorem 2.1 that $\mathbf{G}(\sigma)\mathbf{b} = \widetilde{\mathbf{G}}(\sigma)\mathbf{b}$. The equation $\mathbf{c}^T\mathbf{G}(\sigma) = \mathbf{c}^T\widetilde{\mathbf{G}}(\sigma)$ can be obtained similarly. The proof of part 2 follows from part 3 of Theorem 2.1. $\square$

It should be noted that the conditions $\text{Im}(\mathbf{V}) \subset \text{Im}(\boldsymbol{\Theta}_{r,1})$ and $\text{Im}(\mathbf{W}) \subset \text{Im}(\boldsymbol{\Theta}_{l,2})$ in part 1 of Theorem 6.2 are automatically fulfilled if for given interpolation points $\{\sigma_i\}_{i=1}^r$, $\{\mu_i\}_{i=1}^r$ and tangential directions $\{\mathbf{b}_i\}_{i=1}^r$, $\{\mathbf{c}_i\}_{i=1}^r$, we choose

$$\begin{aligned}
\text{Im}(\mathbf{V}) &= \text{span}\{(\sigma_1\boldsymbol{\mathcal{E}} - \boldsymbol{\mathcal{A}})^I\boldsymbol{\mathcal{B}}\mathbf{b}_1, \ldots, (\sigma_r\boldsymbol{\mathcal{E}} - \boldsymbol{\mathcal{A}})^I\boldsymbol{\mathcal{B}}\mathbf{b}_r\}, \\
\text{Im}(\mathbf{W}) &= \text{span}\{(\mu_1\boldsymbol{\mathcal{E}}^T - \boldsymbol{\mathcal{A}}^T)^I\boldsymbol{\mathcal{C}}^T\mathbf{c}_1, \ldots, (\sigma_r\boldsymbol{\mathcal{E}}^T - \boldsymbol{\mathcal{A}}^T)^I\boldsymbol{\mathcal{C}}^T\mathbf{c}_r\}.
\end{aligned}$$

**6.1. Computational issues related to the reduction of index-2 descriptor systems.** Even though Theorem 6.2 shows how to enforce interpolation for the descriptor system (6.1), the spectral projectors are still implicitly hidden in the definitions of $(\sigma\boldsymbol{\mathcal{E}} - \boldsymbol{\mathcal{A}})^I$ and $(\sigma\boldsymbol{\mathcal{E}}^T - \boldsymbol{\mathcal{A}}^T)^I$. It has been shown in [17] how to compute the matrix-vector product $(\boldsymbol{\mathcal{E}} + \tau\boldsymbol{\mathcal{A}})^I\mathbf{f}$ for a given vector $\mathbf{f}$ without explicitly forming $(\boldsymbol{\mathcal{E}} + \tau\boldsymbol{\mathcal{A}})^I$. This approach can also be used in interpolatory model reduction, where

the quantities of interest are $(\sigma\boldsymbol{\mathcal{E}} - \boldsymbol{\mathcal{A}})^I \boldsymbol{\mathcal{B}} \mathsf{b}$ and $(\mu\boldsymbol{\mathcal{E}}^T - \boldsymbol{\mathcal{A}}^T)^I \boldsymbol{\mathcal{C}}^T \mathsf{c}$. The proof of the following result is analogous to those in [17], and, therefore, it is omitted.

LEMMA 6.3. *Let $s = \sigma, \mu$ be such that $s\boldsymbol{\Theta}_{l,2}^T \mathbf{E}_{11} \boldsymbol{\Theta}_{r,1} - \boldsymbol{\Theta}_{l,2}^T \mathbf{A}_{11} \boldsymbol{\Theta}_{r,1}$ is invertible. Then the vector*

$$\mathbf{v} = (\sigma\boldsymbol{\mathcal{E}} - \boldsymbol{\mathcal{A}})^I \boldsymbol{\mathcal{B}} \mathsf{b} \tag{6.12}$$

*solves*

$$\begin{bmatrix} \sigma\mathbf{E}_{11} - \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{v} \\ \mathbf{z} \end{bmatrix} = \begin{bmatrix} \mathbf{B}_1 \mathsf{b} \\ \mathbf{0} \end{bmatrix}, \tag{6.13}$$

*and the vector*

$$\mathbf{w} = (\mu\boldsymbol{\mathcal{E}}^T - \boldsymbol{\mathcal{A}}^T)^I \boldsymbol{\mathcal{C}}^T \mathsf{c} \tag{6.14}$$

*solves*

$$\begin{bmatrix} \mu\mathbf{E}_{11}^T - \mathbf{A}_{11}^T & \mathbf{A}_{21}^T \\ \mathbf{A}_{12}^T & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ \mathbf{q} \end{bmatrix} = \begin{bmatrix} \mathbf{C}^T \mathsf{c} \\ \mathbf{0} \end{bmatrix}. \tag{6.15}$$

From a computational perspective of implementing Theorem 6.2, the importance of this result is clear. To achieve interpolation, Theorem 6.2 relies on computing the quantities $(\sigma\boldsymbol{\mathcal{E}} - \boldsymbol{\mathcal{A}})^I \boldsymbol{\mathcal{B}} \mathsf{b}$ and $(\sigma\boldsymbol{\mathcal{E}}^T - \boldsymbol{\mathcal{A}}^T)^I \boldsymbol{\mathcal{C}}^T \mathsf{c}$, both of which involve the computation of $\boldsymbol{\Theta}_{l,2}$ and $\boldsymbol{\Theta}_{r,1}$. However, Lemma 6.3 illustrates that the computation of these basis matrices is unnecessary and only the linear systems (6.13) and (6.15) need to be solved. This observation leads to Algorithm 6.1 below for interpolatory model reduction of Stokes-type descriptor systems of index 2.

---

ALGORITHM 6.1. **Interpolatory model reduction for Stokes-type descriptor systems of index 2**

1) *Make an initial selection of the interpolation points $\{\sigma_i\}_{i=1}^r$ and the tangent directions $\{\mathsf{b}_i\}_{i=1}^r$ and $\{\mathsf{c}_i\}_{i=1}^r$.*
2) *For $i = 1, \ldots, r$, solve*

$$\begin{bmatrix} \sigma_i\mathbf{E}_{11} - \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{v}_i \\ \mathbf{z} \end{bmatrix} = \begin{bmatrix} \mathbf{B}_1 \mathsf{b}_i \\ \mathbf{0} \end{bmatrix},$$

$$\begin{bmatrix} \sigma_i\mathbf{E}_{11}^T - \mathbf{A}_{11}^T & \mathbf{A}_{21}^T \\ \mathbf{A}_{12}^T & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{w}_i \\ \mathbf{q} \end{bmatrix} = \begin{bmatrix} \mathbf{C}^T \mathsf{c}_i \\ \mathbf{0} \end{bmatrix}.$$

3) $\mathbf{V} = [\mathbf{v}_1, \ldots, \mathbf{v}_r], \quad \mathbf{W} = [\mathbf{w}_1, \ldots, \mathbf{w}_r]$.
4) $\widetilde{\mathbf{E}} = \mathbf{W}^T \mathbf{E}_{11} \mathbf{V}, \ \widetilde{\mathbf{A}} = \mathbf{W}^T \mathbf{A}_{11} \mathbf{V}, \ \widetilde{\mathbf{B}} = \mathbf{W}^T \mathbf{B}_1, \ \widetilde{\mathbf{C}} = \mathbf{C}\mathbf{V}, \ \widetilde{\mathbf{D}} = \boldsymbol{\mathfrak{D}}$.

---

Once a computationally effective bi-tangential Hermite interpolation framework is established for index-2 descriptor systems, extending it to optimal $\mathcal{H}_2$ model reduction via IRKA is straightforward and given in Algorithm 6.2. It follows from the structure of this algorithm that upon convergence the reduced model $\widetilde{\mathbf{G}}(s) = \widetilde{\mathbf{C}}(s\widetilde{\mathbf{E}} - \widetilde{\mathbf{A}})^{-1}\widetilde{\mathbf{B}} + \widetilde{\mathbf{D}}$ satisfies the first-order conditions for $\mathcal{H}_2$ optimality.

REMARK 6.2. *As shown in [17], the general case $\mathbf{B}_2 \neq \mathbf{0}$ can be handled similar to the case $\mathbf{B}_2 = \mathbf{0}$. First note that the state $\mathbf{x}_1(t)$ can be decomposed as*

---

ALGORITHM 6.2. **IRKA for Stokes-type descriptor system of index 2**

1) *Make an initial shift selection $\{\sigma_i\}_{i=1}^r$ and initial tangent directions $\{b_i\}_{i=1}^r$ and $\{c_i\}_{i=1}^r$.*
2) *Apply Algorithm 6.1 to obtain $\widetilde{\mathbf{E}}$, $\widetilde{\mathbf{A}}$, $\widetilde{\mathbf{B}}$, $\widetilde{\mathbf{C}}$ and $\widetilde{\mathbf{D}}$.*
3) *while (not converged)*
   a) *Compute $\mathbf{Y}^*\widetilde{\mathbf{A}}\mathbf{Z} = diag(\lambda_1, \ldots, \lambda_r)$ and $\mathbf{Y}^*\widetilde{\mathbf{E}}\mathbf{Z} = \mathbf{I}$, where the columns of $\mathbf{Z} = [\mathbf{z}_1, \ldots, \mathbf{z}_r]$ and $\mathbf{Y} = [\mathbf{y}_1, \ldots, \mathbf{y}_r]$ are, respectively, the right and left eigenvectors of $\lambda\widetilde{\mathbf{E}} - \widetilde{\mathbf{A}}$.*
   b) *$\sigma_i \leftarrow -\lambda_i$, $b_i^T \leftarrow \mathbf{y}_i^*\widetilde{\mathbf{B}}$ and $c_i \leftarrow \widetilde{\mathbf{C}}\mathbf{z}_i$ for $i = 1, \ldots, r$.*
   c) *Apply Algorithm 6.1 to obtain $\widetilde{\mathbf{E}}$, $\widetilde{\mathbf{A}}$, $\widetilde{\mathbf{B}}$, $\widetilde{\mathbf{C}}$ and $\widetilde{\mathbf{D}}$.*
   *end while*

---

$\mathbf{x}_1(t) = \mathbf{x}_0(t) + \mathbf{x}_g(t)$, where $\mathbf{x}_g(t) = -\mathbf{E}_{11}^{-1}\mathbf{A}_{12}(\mathbf{A}_{21}\mathbf{E}_{11}^{-1}\mathbf{A}_{12})^{-1}\mathbf{B}_2\mathbf{u}(t)$ and $\mathbf{x}_0(t)$ satisfies $\mathbf{A}_{21}\mathbf{x}_0(t) = 0$. After some algebraic manipulations, this leads to

$$\begin{aligned}\mathbf{\Pi}_l\mathbf{E}_{11}\mathbf{\Pi}_r\dot{\mathbf{x}}_0(t) &= \mathbf{\Pi}_l\mathbf{A}_{11}\mathbf{\Pi}_r\mathbf{x}_0(t) + \mathbf{\Pi}_l\mathbf{B}\mathbf{u}(t),\\ \mathbf{y}(t) &= \mathbf{C}\mathbf{\Pi}_r\mathbf{x}_0(t) + \boldsymbol{\mathcal{D}}\mathbf{u}(t) - \mathbf{C}_2(\mathbf{A}_{21}\mathbf{E}_{11}^{-1}\mathbf{A}_{12})^{-1}\mathbf{B}_2\dot{\mathbf{u}}(t),\end{aligned} \tag{6.16}$$

*where*

$$\mathbf{C} = \mathbf{C}_1 - \mathbf{C}_2(\mathbf{A}_{21}\mathbf{E}_{11}^{-1}\mathbf{A}_{12})^{-1}\mathbf{A}_{21}\mathbf{E}_{11}^{-1}\mathbf{A}_{11}, \tag{6.17}$$

$$\mathbf{B} = \mathbf{B}_1 - \mathbf{A}_{11}\mathbf{E}_{11}^{-1}\mathbf{A}_{12}(\mathbf{A}_{21}\mathbf{E}_{11}^{-1}\mathbf{A}_{12})^{-1}\mathbf{B}_2, \tag{6.18}$$

$$\boldsymbol{\mathcal{D}} = \mathbf{D} - \mathbf{C}_2(\mathbf{A}_{21}\mathbf{E}_{11}^{-1}\mathbf{A}_{12})^{-1}\mathbf{A}_{21}\mathbf{E}_{11}^{-1}\mathbf{B}_1. \tag{6.19}$$

*Therefore, the $\mathbf{B}_2 \neq 0$ case extends to the interpolation framework as well by defining*

$$\widehat{\boldsymbol{\mathcal{E}}} = \mathbf{\Pi}_l\mathbf{E}_{11}\mathbf{\Pi}_r, \quad \widehat{\boldsymbol{\mathcal{A}}} = \mathbf{\Pi}_l\mathbf{A}_{11}\mathbf{\Pi}_r, \quad \widehat{\boldsymbol{\mathcal{B}}} = \mathbf{\Pi}_l\mathbf{B}, \quad \widehat{\boldsymbol{\mathcal{C}}} = \mathbf{C}\mathbf{\Pi}_r$$

*and applying Theorem 6.2 with $\widehat{\boldsymbol{\mathcal{E}}}$, $\widehat{\boldsymbol{\mathcal{A}}}$, $\widehat{\boldsymbol{\mathcal{B}}}$, $\widehat{\boldsymbol{\mathcal{C}}}$ and $\widehat{\boldsymbol{\mathcal{D}}} = \boldsymbol{\mathcal{D}} - s\mathbf{C}_2(\mathbf{A}_{21}\mathbf{E}_{11}^{-1}\mathbf{A}_{12})^{-1}\mathbf{B}_2$ instead of $\boldsymbol{\mathcal{E}}$, $\boldsymbol{\mathcal{A}}$, $\boldsymbol{\mathcal{B}}$, $\boldsymbol{\mathcal{C}}$ and $\boldsymbol{\mathcal{D}}$.*

**6.2. Numerical results for Oseen equations.** The model borrowed from [17] is obtained by discretizing the Oseen equations and describe the flow of a viscous and incompressible fluid in a domain $\Omega \in \mathbb{R}^2$ representing a channel with a backward facing step. A spatial discretization using the finite element method leads to the index-2 descriptor system (6.1) with $\mathbf{E}_{11}, \mathbf{A}_{11} \in \mathbb{R}^{5520 \times 5520}$, $\mathbf{A}_{12}, \mathbf{A}_{21}^T \in \mathbb{R}^{5520 \times 761}$, $\mathbf{B}_1 \in \mathbb{R}^{5520 \times 6}$, $\mathbf{B}_2 \in \mathbb{R}^{761 \times 6}$, $\mathbf{C}_1 \in \mathbb{R}^{2 \times 5520}$, $\mathbf{C}_2 \in \mathbb{R}^{2 \times 761}$, $\mathbf{D} = 0$, see [17] for more details on the model. Note that $\mathbf{B}_2 \neq 0$ and the transfer function grows unbounded around $s = \infty$.

We approximate this system by a model of order $r = 20$ using the balanced truncation method as described in [17] and the $\mathcal{H}_2$ optimal model reduction method given in Algorithm 6.2. The amplitude Bode plots of the full model and two reduced-order models depicted in Figure 6.1 clearly illustrate that interpolation-based Algorithm 6.2 leads to a high-fidelity reduced model replicating the full-order transfer function with almost no loss of accuracy and matching the performance of the balanced truncation method. The accuracy of this interpolation-based method is due to the fact that we do not choose the interpolation points in an *ad hoc* fashion; instead Algorithm 6.2 iteratively leads to $\mathcal{H}_2$ optimal interpolation points. As the difficulty in computing $\mathcal{H}_2$ norm of the error is clear, we approximately compute the relative $\mathcal{H}_\infty$-error $\frac{\|\mathbf{G}_{sp} - \tilde{\mathbf{G}}_{sp}\|_{\mathcal{H}_\infty}}{\|\mathbf{G}_{sp}\|_{\mathcal{H}_\infty}}$ for both reduced-order models by sampling the imaginary axis.

These errors for the balanced truncation method and Algorithm 6.2 are, respectively, $3.3284 \times 10^{-6}$ and $8.9663 \times 10^{-6}$. Both reduced-order models are highly accurate. It is expected that the $\mathcal{H}_\infty$-error in balanced truncation will be smaller than that in IRKA. While our method tries to minimize the $\mathcal{H}_2$-norm, the balanced truncation method is tailored towards reducing the $\mathcal{H}_\infty$-norm. Indeed, these numbers are further signs for the success of the interpolatory-based model reduction method as it produces a very accurate model, almost matching the accuracy of the balanced truncation approach. These observations are similar to those on IRKA whose $\mathcal{H}_\infty$-norm behavior was close to or even better in some cases than that of balanced truncation [1, 15].



Fig. 6.1. *Oseen equation: amplitude Bode plots of the full and reduced models*

To further illustrate the accuracy in the reduced-order model computed by Algorithm 6.2, we display the time domain response plots resulting from two different input selections. In the left pane of Figure 6.2, we plot the outputs for the input selections $\mathbf{u}_i(t) = \sin(6it)$ for $i = 1, \ldots, 6$ (recall that the system has 6 inputs). The figure illustrate a perfect match between the outputs of the full and reduced-order systems. Error in the outputs for the same input selection is given in the right pane of Figure 6.2. Note the difference in the scale of the error plot compared to the actual output; the error is four orders of magnitude smaller. We repeat the same experiments with $\mathbf{u}_i(t) = \sin(it)$ for $i = 1, \ldots, 6$ and reach the same conclusions as shown in Figure 6.3.

**7. Conclusions.** For interpolatory model reduction of descriptor systems, we have introduced subspace conditions that not only guarantee interpolation conditions but also automatically enforce matching the polynomial part of the transfer function, thus preventing the error grow unbounded. We have also extended the optimal $\mathcal{H}_2$ interpolation point selection strategy to descriptor systems. For the index-1 and index-2 descriptor systems, we have shown how to construct the reduced-order models without computing the deflating subspaces corresponding to the finite and infinite eigenvalues explicitly. Several numerical examples have supported the theoretical discussion.
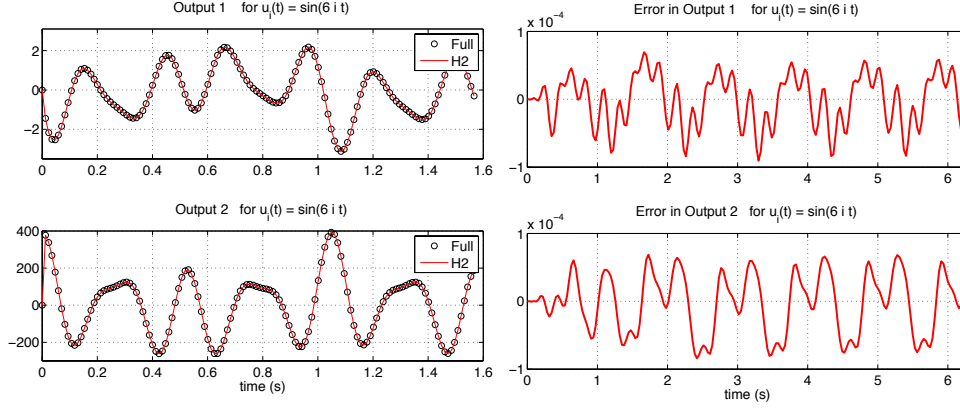
FIG. 6.2. *Oseen equation: (left) time domain response for* $\mathbf{u}_i(t) = \sin(6it)$*; (right) error in time domain response for* $\mathbf{u}_i(t) = \sin(6it)$.
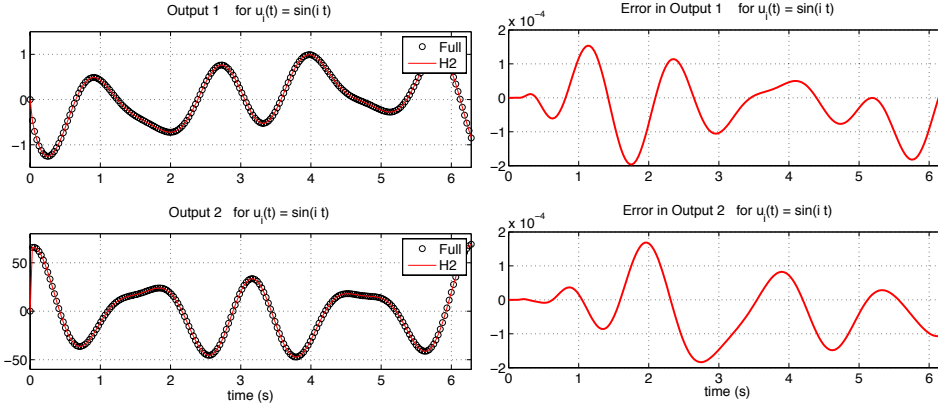


FIG. 6.3. *Oseen equation: (left) time domain response for* $\mathbf{u}_i(t) = \sin(it)$*; (right) error in time domain response for* $\mathbf{u}_i(t) = \sin(it)$.

## REFERENCES

[1] A.C. Antoulas, C.A. Beattie, and S. Gugercin. Interpolatory model reduction of large-scale dynamical systems. In J. Mohammadpour and K. Grigoriadis, editors, *Efficient Modeling and Control of Large-Scale Systems*, pages 3–58. Springer-Verlag, 2010.

[2] C. Beattie and S. Gugercin. Interpolatory projection methods for structure-preserving model reduction. *Systems Control Lett.*, 58(3):225–232, 2009.

[3] C.A. Beattie and S. Gugercin. Krylov-based minimization for optimal $\mathcal{H}_2$ model reduction. *Proceedings of the 46th IEEE Conference on Decision and Control*, pages 4385–4390, 2007.

[4] C.A. Beattie and S. Gugercin. A trust region method for optimal $\mathcal{H}_2$ model reduction. *Proceeding of the 48th IEEE Conference on Decision and Control*, 2009.

[5] P. Benner and V.I. Sokolov. Partial realization of descriptor systems. *Systems Control Lett.*, 55(11):929–938, 2006.

[6] A. Bunse-Gerstner, D. Kubalinska, G. Vossen, and D. Wilczek. $\mathcal{H}_2$-optimal model reduction for large scale discrete dynamical MIMO systems. *J. Comput. Appl. Math.*, 233(5):1202–1216, 2010.

[7] C. De Villemagne and R.E. Skelton. Model reductions using a projection formulation. *Intern. J. Control*, 46(6):2141–2169, 1987.

[8] P. Feldmann and R.W. Freund. Efficient linear circuit analysis by Padé approximation via the Lanczos process. *IEEE Trans. Computer-Aided Design Integr. Circuits Syst.*, 14(5):639–

649, 1995.

[9] K. Gallivan, A. Vandendorpe, and P. Van Dooren. Model reduction of MIMO systems via tangential interpolation. *SIAM J. Matrix Anal. Appl.*, 26(2):328–349, 2005.

[10] E. Grimme. *Krylov Projection Methods for Model Reduction*. PhD thesis, University of Illinois, Urbana-Champaign, 1997.

[11] S. Gugercin. *Projection methods for model reduction of large-scale dynamical systems*. PhD thesis, Rice University, 2002.

[12] S. Gugercin. An iterative rational Krylov algorithm (IRKA) for optimal $\mathcal{H}_2$ model reduction. In *Householder Symposium XVI*, Seven Springs Mountain Resort, PA, USA, May 2005.

[13] S. Gugercin and A.C. Antoulas. An $\mathcal{H}_2$ error expression for the Lanczos procedure. In *Proceedings of the 42nd IEEE Conference on Decision and Control*, 2003.

[14] S. Gugercin, A.C. Antoulas, and C.A. Beattie. A rational Krylov iteration for optimal $\mathcal{H}_2$ model reduction. In *Proceedings of 17th International Symposium on Mathematical Theory of Networks and Systems* (July 24-28, 2006, Kyoto, Japan), 2006.

[15] S. Gugercin, A.C. Antoulas, and C.A. Beattie. $\mathcal{H}_2$ model reduction for large-scale linear dynamical systems. *SIAM J. Matrix Anal. Appl.*, 30(2):609–638, 2008.

[16] Y. Halevi. Frequency weighted model reduction via optimal projection. *IEEE Trans. Automat. Control*, 37(10):1537–1542, 1992.

[17] M. Heinkenschloss, D.C. Sorensen, and K. Sun. Balanced truncation model reduction for a class of descriptor systems with application to the oseen equations. *SIAM J. Sci. Comput.*, 30(2):1038–1063, 2008.

[18] D. Hyland and D. Bernstein. The optimal projection equations for model reduction and the relationships among the methods of Wilson, Skelton, and Moore. *IEEE Trans. Automat. Control*, 30(12):1201–1211, 1985.

[19] D. Kubalinska, A. Bunse-Gerstner, G. Vossen, and D. Wilczek. $\mathcal{H}_2$-optimal interpolation based model reduction for large-scale systems. In *Proceedings of the* 16[th] *International Conference on System Science*, Poland, 2007.

[20] P. Kunkel and V. Mehrmann. *Differential-Algebraic Equations. Analysis and Numerical Solution*. EMS Publishing House, Zürich, Switzerland, 2006.

[21] G. Lassaux and K. Willcox. Model reduction of an actively controlled supersonic diffuser. In P. Benner, V. Mehrmann, and D. C. Sorensen, editors, *Dimension Reduction of Large-Scale Systems*, volume 45 of *Lecture Notes in Computational Science and Engineering*, pages 357–361. Springer-Verlag, Berlin, Heidelberg, Germany, 2005.

[22] A.J. Mayo and A.C. Antoulas. A framework for the solution of the generalized realization problem. *Linear Algebra Appl.*, 425(2-3):634–662, 2007.

[23] L. Meier III and D. Luenberger. Approximation of linear constant systems. *IEEE Trans. Automat. Control*, 12(5):585–588, 1967.

[24] A. Ruhe. Rational Krylov algorithms for nonsymmetric eigenvalue problems. II: Matrix pair. *Linear Algebra Appl.*, 197-198:282–295, 1994.

[25] J.T. Spanos, M.H. Milman, and D.L. Mingori. A new algorithm for $L^2$ optimal model reduction. *Automatica*, 28(5):897–909, 1992.

[26] T. Stykel. Gramian-based model reduction for descriptor systems. *Math. Control Signals Syst.*, 16(4):297–319, 2004.

[27] T. Stykel. Low-rank iterative methods for projected generalized Lyapunov equations. *Electron. Trans. Numer. Anal.*, 30:187–202, 2008.

[28] P. Van Dooren, K.A. Gallivan, and P.A. Absil. $\mathcal{H}_2$-optimal model reduction of MIMO systems. *Appl. Math. Lett.*, 21(12):1267–1273, 2008.

[29] D.A. Wilson. Optimum solution of model-reduction problem. *Proc. IEE*, 117(6):1161–1165, 1970.

[30] W.Y. Yan and J. Lam. An approximate approach to $H^2$ optimal model reduction. *IEEE Trans. Automat. Control*, 44(7):1341–1358, 1999.

[31] A. Yousouff, D.A. Wagie, and R.E. Skelton. Linear system approximation via covariance equivalent realizations. *J. Math. Anal. Appl.*, 196:91–115, 1985.

[32] A. Yousuff and R.E. Skelton. Covariance equivalent realizations with applications to model reduction of large-scale systems. In C.T. Leondes, editor, *Control and Dynamic Systems*, volume 22, pages 273–348. Academic Press, 1985.

[33] D. Zigic, L.T. Watson, and C. Beattie. Contragredient transformations applied to the optimal projection equations. *Linear Algebra Appl.*, 188:665–676, 1993.